

# Data Mining

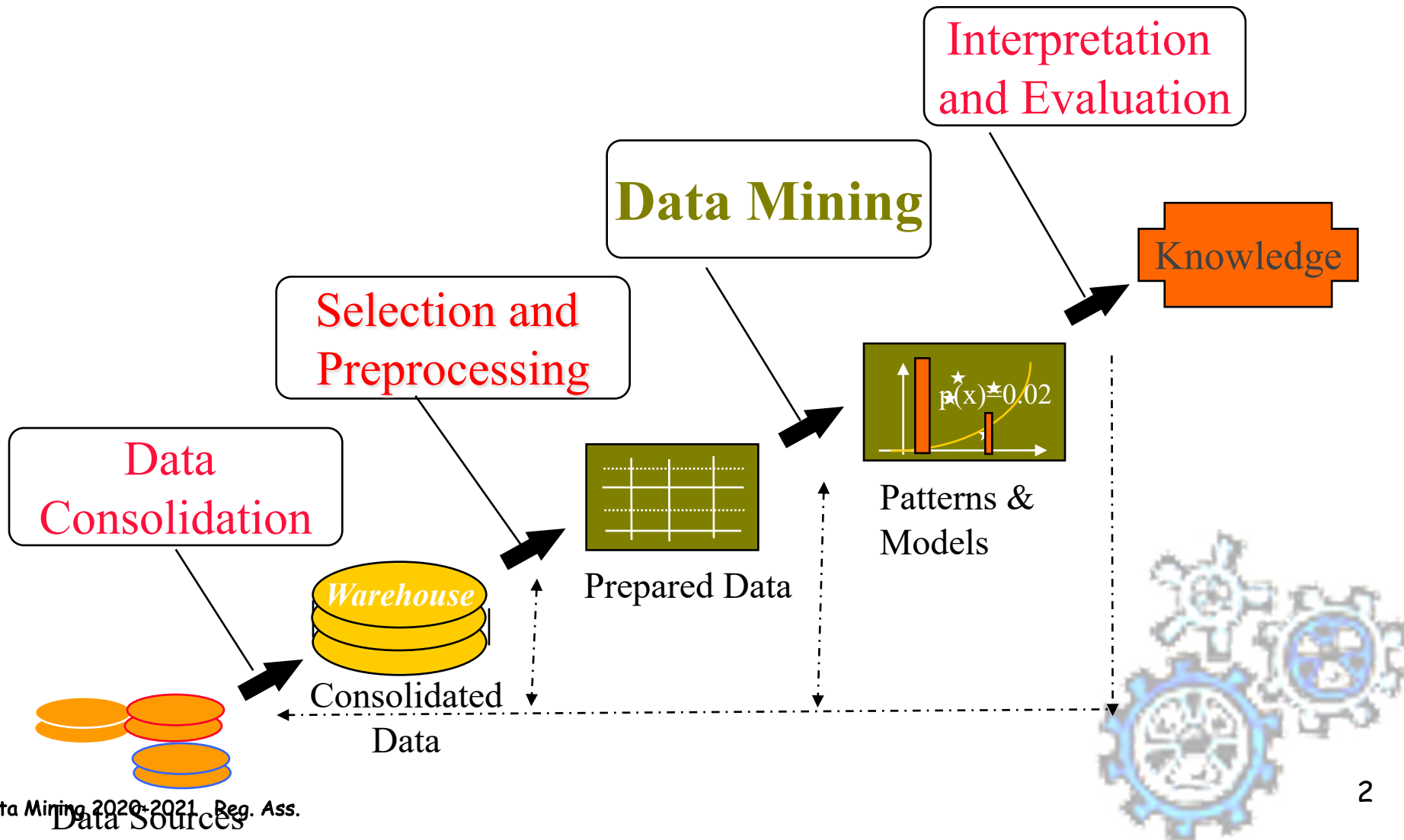
## Knowledge Discovery in Databases

**Dino Pedreschi, Mirco Nanni**  
**Pisa KDD Lab, ISTI-CNR & Univ. Pisa**



Data Mining 2020-2021  
Reg. Ass.

# KDD Process



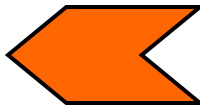


# Association rules and market basket analysis



# Association rules - module outline

1. What are association rules (AR) and what are they used for:
  1. The paradigmatic application: Market Basket Analysis
  2. The single dimensional AR (intra-attribute)
2. How to compute AR
  1. Basic Apriori Algorithm and its optimizations
  2. Multi-Dimension AR (inter-attribute)
  3. Quantitative AR
  4. Constrained AR
3. How to reason on AR and how to evaluate their quality
  1. Multiple-level AR
  2. Interestingness
  3. Correlation vs. Association



# Market Basket Analysis: the context

Customer buying habits by finding associations and correlations between the different items that customers place in their “shopping basket”

Milk, eggs, sugar,  
bread



Customer1

Milk, eggs, cereal, bread



Customer2

Eggs, sugar



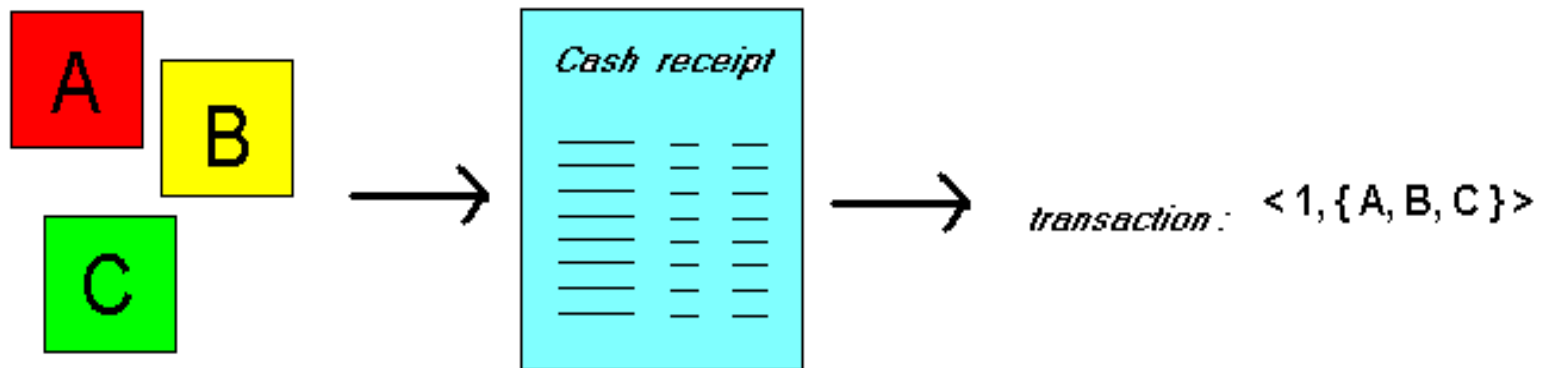
Customer3



# Market Basket Analysis: the context

Given: a database of customer **transactions**, where each transaction is a **set of items**

- Find groups of items which are **frequently purchased together**



# Goal of MBA

- Extract information on purchasing behavior
- Actionable information: can suggest
  - new store layouts
  - new product assortments
  - which products to put on promotion
- MBA applicable whenever a customer purchases multiple things in proximity
  - credit cards
  - services of telecommunication companies
  - banking services
  - medical treatments



# MBA: applicable to many other contexts

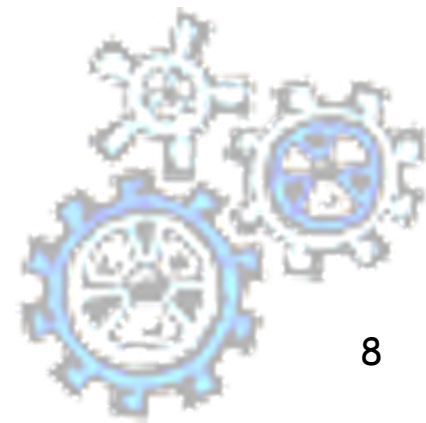
## Telecommunication:

Each customer is a transaction containing the set of customer's phone calls

## Atmospheric phenomena:

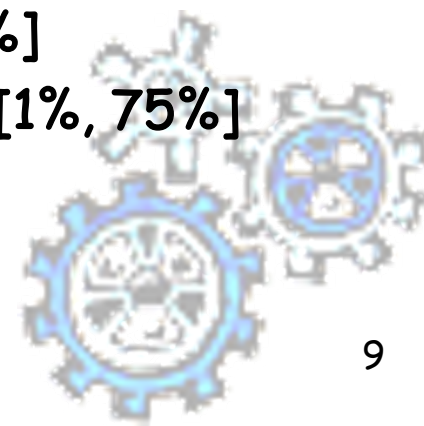
Each time interval (e.g. a day) is a transaction containing the set of observed event (rains, wind, etc.)

Etc.



# Association Rules

- Express how product/services relate to each other, and tend to group together
- “if a customer purchases three-way calling, then will also purchase call-waiting”
- simple to understand
- actionable information: bundle three-way calling and call-waiting in a single package
- Examples.
  - Rule form: “Body → Head [support, confidence]”.
  - $\text{buys}(x, \text{"diapers"}) \rightarrow \text{buys}(x, \text{"beers"}) [0.5\%, 60\%]$
  - $\text{major}(x, \text{"CS"}) \wedge \text{takes}(x, \text{"DB"}) \rightarrow \text{grade}(x, \text{"A"}) [1\%, 75\%]$



# Useful, trivial, unexplicable

- **Useful:** “On Thursdays, grocery store consumers often purchase diapers and beer together”.
- **Trivial:** “Customers who purchase maintenance agreements are very likely to purchase large appliances”.
- **Unexplicable:** “When a new hardware store opens, one of the most sold items is toilet rings.”





# Association Rules Road Map

- Single dimension vs. multiple dimensional AR
  - E.g., association on items bought vs. linking on different attributes.
  - Intra-Attribute vs. Inter-Attribute
- Qualitative vs. quantitative AR
  - Association on categorical vs. numerical attributes
- Simple vs. constraint-based AR
  - E.g., small sales ( $\text{sum} < 100$ ) trigger big buys ( $\text{sum} > 1,000$ )?
- Single level vs. multiple-level AR
  - E.g., what **brands** of beers are associated with what **brands** of diapers?
- Association vs. correlation analysis.
  - Association does not necessarily imply correlation.



# Association rules - module outline

- What are association rules (AR) and what are they used for:
  - The paradigmatic application: Market Basket Analysis
  - The single dimensional AR (intra-attribute)
- How to compute AR
  - Basic Apriori Algorithm and its optimizations
  - Multi-Dimension AR (inter-attribute)
  - Quantitative AR
  - Constrained AR
- How to reason on AR and how to evaluate their quality
  - Multiple-level AR
  - Interestingness
  - Correlation vs. Association



# Data Mining

## Association Analysis: Basic Concepts and Algorithms

---

### Lecture Notes for Chapter 6

Introduction to Data Mining  
by  
Tan, Steinbach, Kumar

# Association Rule Mining

- Given a set of transactions, find rules that will predict the occurrence of an item based on the occurrences of other items in the transaction

## Market-Basket transactions

<i>TID</i>	<i>Items</i>
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

## Example of Association Rules

$\{\text{Diaper}\} \rightarrow \{\text{Beer}\},$   
 $\{\text{Milk, Bread}\} \rightarrow \{\text{Eggs, Coke}\},$   
 $\{\text{Beer, Bread}\} \rightarrow \{\text{Milk}\},$

Implication means co-occurrence,  
not causality!



# Definition: Frequent Itemset

## ■ Itemset

- A collection of one or more items

- ✓ Example: {Milk, Bread, Diaper}

## ■ k-itemset

- ✓ An itemset that contains k items

## ■ Support count ( $\sigma$ )

- Frequency of occurrence of an itemset
- E.g.  $\sigma(\{\text{Milk, Bread, Diaper}\}) = 2$
- $\sigma(X) = |\{t_i | X \text{ contained in } t_i \text{ and } t_i \text{ is a transaction}\}|$

## ■ Support

- Fraction of transactions that contain an itemset
- E.g.  $s(\{\text{Milk, Bread, Diaper}\}) = 2/5$

## ■ Frequent Itemset

- An itemset whose support is greater than or equal to a *minsup* threshold

<i>TID</i>	<i>Items</i>
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke



# Definition: Association Rule

## ■ Association Rule

- An implication expression of the form  $X \rightarrow Y$ , where  $X$  and  $Y$  are itemsets
- Example:  
 $\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$

<i>TID</i>	<i>Items</i>
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

## ■ Rule Evaluation Metrics

- **Support (s)**
  - ✓ Fraction of transactions that contain both  $X$  and  $Y$
- **Confidence (c)**
  - ✓ Measures how often items in  $Y$  appear in transactions that contain  $X$

Example:

$\{\text{Milk, Diaper}\} \Rightarrow \text{Beer}$

$$s = \frac{\sigma(\text{Milk, Diaper, Beer})}{|T|} = \frac{2}{5} = 0.4$$

$$c = \frac{\sigma(\text{Milk, Diaper, Beer})}{\sigma(\text{Milk, Diaper})} = \frac{2}{3} = 0.67$$

# Association Rule Mining Task

- Given a set of transactions  $T$ , the goal of association rule mining is to find all rules having
  - support  $\geq$  *minsup* threshold
  - confidence  $\geq$  *minconf* threshold
- Brute-force approach:
  - List all possible association rules
  - Compute the support and confidence for each rule
  - Prune rules that fail the *minsup* and *minconf* thresholds

⇒ **Computationally prohibitive!**



# Mining Association Rules

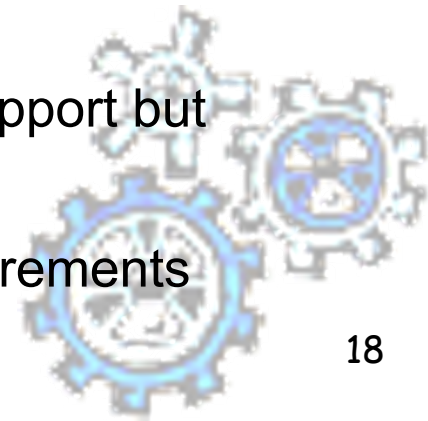
<i>TID</i>	<i>Items</i>
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

## Example of Rules:

$\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$  ( $s=0.4$ ,  $c=0.67$ )  
 $\{\text{Milk, Beer}\} \rightarrow \{\text{Diaper}\}$  ( $s=0.4$ ,  $c=1.0$ )  
 $\{\text{Diaper, Beer}\} \rightarrow \{\text{Milk}\}$  ( $s=0.4$ ,  $c=0.67$ )  
 $\{\text{Beer}\} \rightarrow \{\text{Milk, Diaper}\}$  ( $s=0.4$ ,  $c=0.67$ )  
 $\{\text{Diaper}\} \rightarrow \{\text{Milk, Beer}\}$  ( $s=0.4$ ,  $c=0.5$ )  
 $\{\text{Milk}\} \rightarrow \{\text{Diaper, Beer}\}$  ( $s=0.4$ ,  $c=0.5$ )

## Observations:

- All the above rules are binary partitions of the same itemset:  
 $\{\text{Milk, Diaper, Beer}\}$
- Rules originating from the same itemset have identical support but can have different confidence
- Thus, we may decouple the support and confidence requirements





# Mining Association Rules

## ■ Two-step approach:

### 1. Frequent Itemset Generation

- Generate all itemsets whose support  $\geq$  minsup

### 2. Rule Generation

- Generate high confidence rules from each frequent itemset, where each rule is a binary partitioning of a frequent itemset

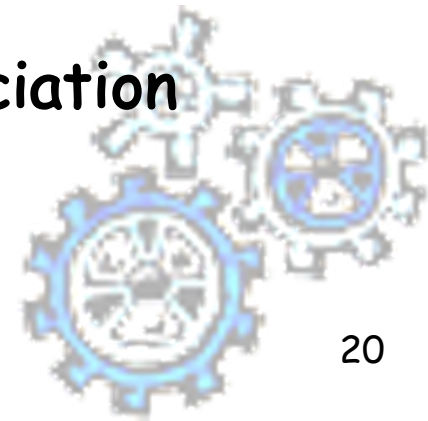
## ■ Frequent itemset generation is still computationally expensive



# Basic Apriori Algorithm

## Problem Decomposition

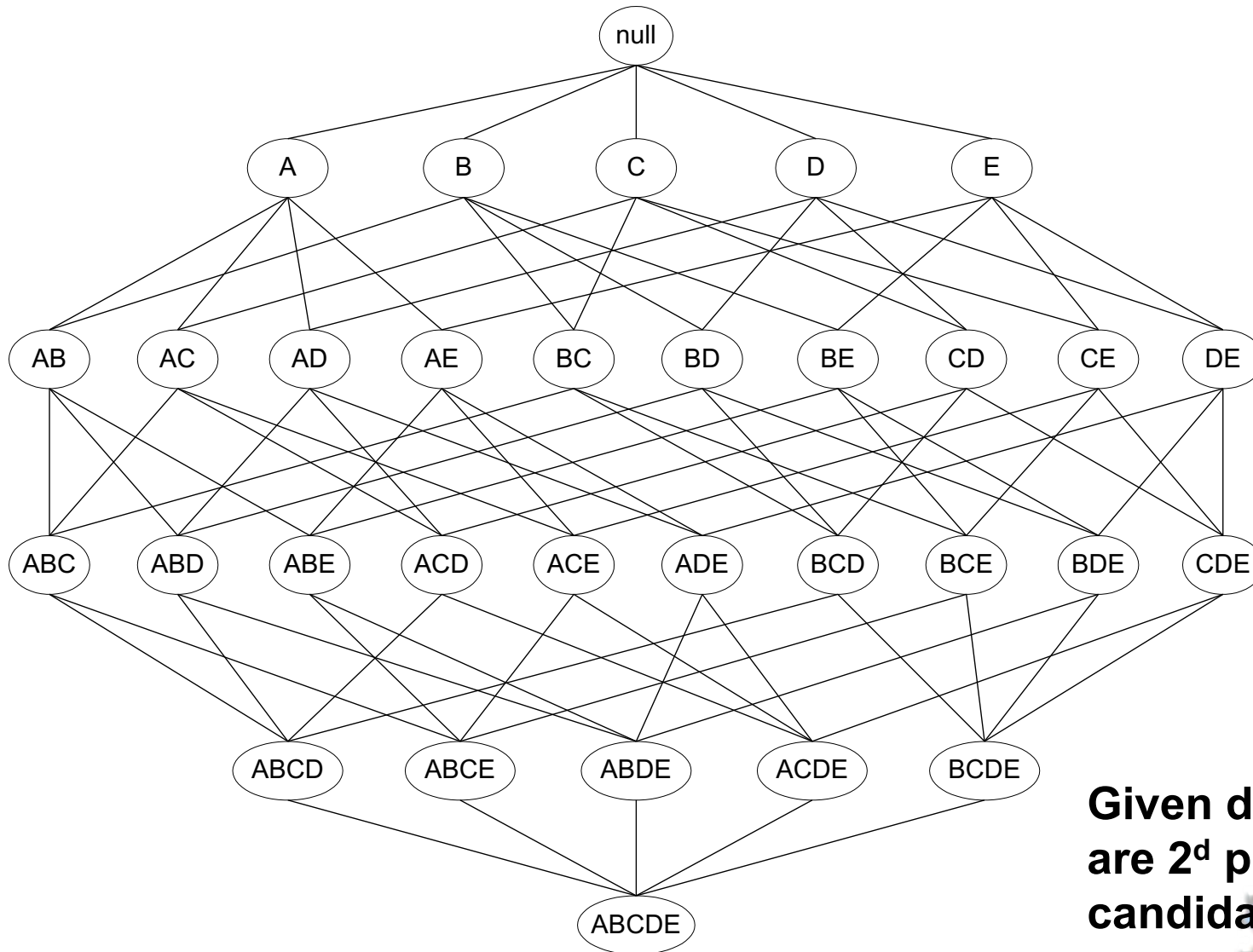
- ① Find the *frequent itemsets*: the sets of items that satisfy the support constraint
  - ◆ A subset of a frequent itemset is also a frequent itemset, i.e., if  $\{A, B\}$  is a frequent itemset, both  $\{A\}$  and  $\{B\}$  should be a frequent itemset
  - ◆ Iteratively find frequent itemsets with cardinality from 1 to  $k$  ( $k$ -itemset)
- ② Use the frequent itemsets to generate association rules.



# Frequent Itemset Mining Problem

- $I = \{x_1, \dots, x_n\}$  set of distinct literals (called **items**)
- $X \subseteq I, X \neq \emptyset, |X| = k, X$  is called **k-itemset**
- A **transaction** is a couple  $\langle tID, X \rangle$  where  $X$  is an itemset
- A **transaction database**  $TDB$  is a set of transactions
- An itemset  $X$  is **contained** in a transaction  $\langle tID, Y \rangle$  if  $X \subseteq Y$
- Given a  $TDB$  the subset of transactions of  $TDB$  in which  $X$  is contained is named  $TDB[X]$ .
- The **support(COUNT)** of an itemset  $X$ , written  $supp_{TDB}(X)$  is the cardinality of  $TDB[X]$ .
- **The support(relative)** of an itemset  $X$ , written  $supp(X)$  is the cardinality of  $TDB[X]$  / cardinality of  $TDB$ .
- Given a user-defined **min\_sup** threshold an itemset  $X$  is **frequent** in  $TDB$  if its support is no less than **min\_sup**.
- Given a user-defined **min\_sup** and a transaction database  $TDB$ , the **Frequent Itemset Mining Problem** requires to compute all frequent itemsets in  $TDB$  w.r.t **min\_sup**.

# Frequent Itemset Generation

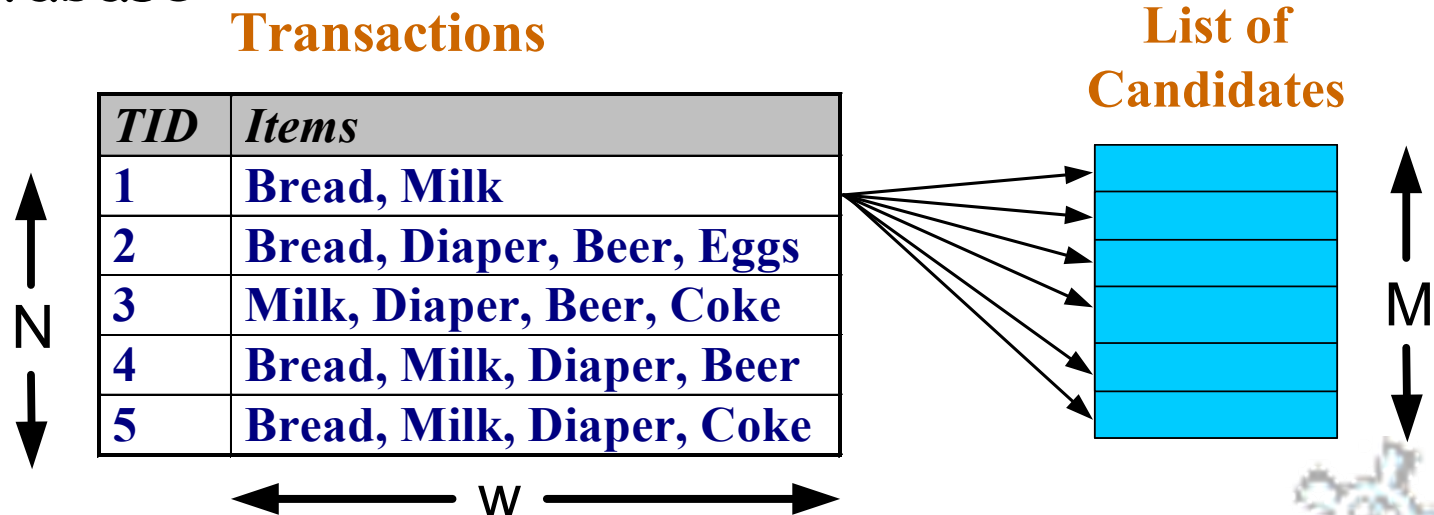


**Given  $d$  items, there are  $2^d$  possible candidate itemsets**

# Frequent Itemset Generation

## ■ Brute-force approach:

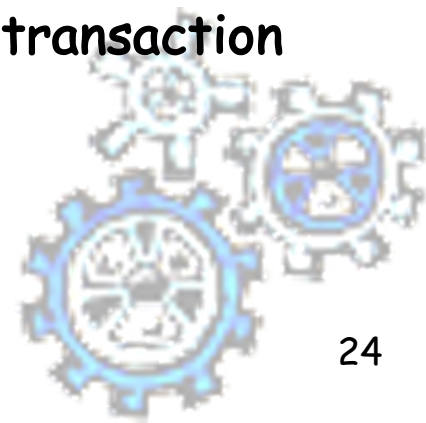
- Each itemset in the lattice is a **candidate** frequent itemset
- Count the support of each candidate by scanning the database



- Match each transaction against every candidate
- Complexity  $\sim O(NMw) \Rightarrow$  **Expensive since  $M = 2^d$  !!!**

# Frequent Itemset Generation Strategies

- Reduce the **number of candidates** ( $M$ )
  - Complete search:  $M=2^d$
  - Use pruning techniques to reduce  $M$
- Reduce the **number of transactions** ( $N$ )
  - Reduce size of  $N$  as the size of itemset increases
  - Used by DHP and vertical-based mining algorithms
- Reduce the **number of comparisons** ( $NM$ )
  - Use efficient data structures to store the candidates or transactions
  - No need to match every candidate against every transaction



# Reducing Number of Candidates

## ■ Apriori principle:

- If an itemset is frequent, then all of its subsets must also be frequent

## ■ Apriori principle holds due to the following property of the support measure:

$$\forall X, Y : (X \subseteq Y) \Rightarrow s(X) \geq s(Y)$$

- Support of an itemset never exceeds the support of its subsets
- This is known as the **anti-monotone** property of support



# The Apriori property

- If  $B$  is frequent and  $A \subseteq B$  then  $A$  is also frequent
  - Each transaction which contains  $B$  contains also  $A$ , which implies  $\text{supp.}(A) \geq \text{supp.}(B)$
- **Consequence:** if  $A$  is not frequent, then it is not necessary to generate the itemsets which include  $A$ .
- **Example:**

•  $\langle 1, \{a, b\} \rangle$        $\langle 2, \{a\} \rangle$   
•  $\langle 3, \{a, b, c\} \rangle$        $\langle 4, \{a, b, d\} \rangle$

with minimum support = 30%.

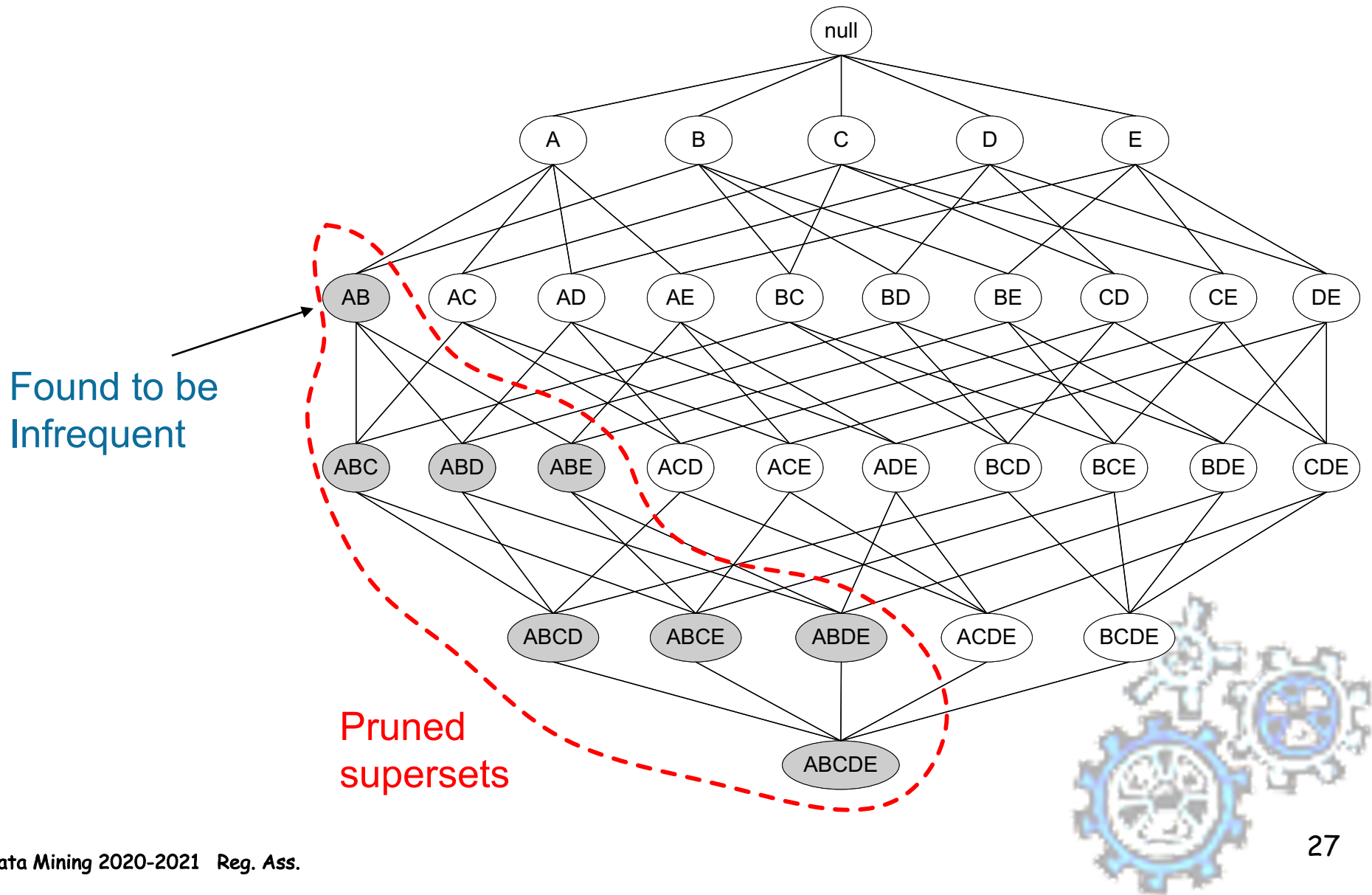
The itemset  $\{c\}$  is not frequent so is not necessary to check for:

$\{c, a\}, \{c, b\}, \{c, d\}, \{c, a, b\}, \{c, a, d\}, \{c, b, d\}$





# Illustrating Apriori Principle



# Illustrating Apriori Principle

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Items (1-itemsets)



Itemset	Count
{Bread,Milk}	3
{Bread,Beer}	2
{Bread,Diaper}	3
{Milk,Beer}	2
{Milk,Diaper}	3
{Beer,Diaper}	3

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

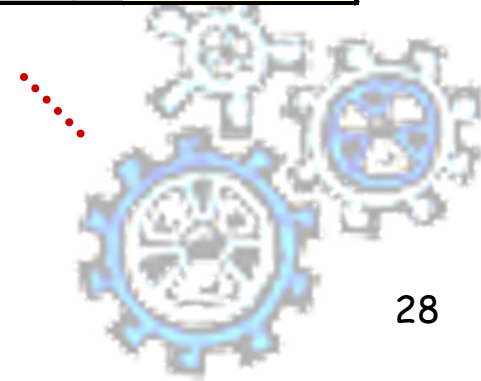
Minimum Support = 3

If every subset is considered,  
 ${}^6C_1 + {}^6C_2 + {}^6C_3 = 41$   
With support-based pruning,  
 $6 + 6 + 1 = 13$

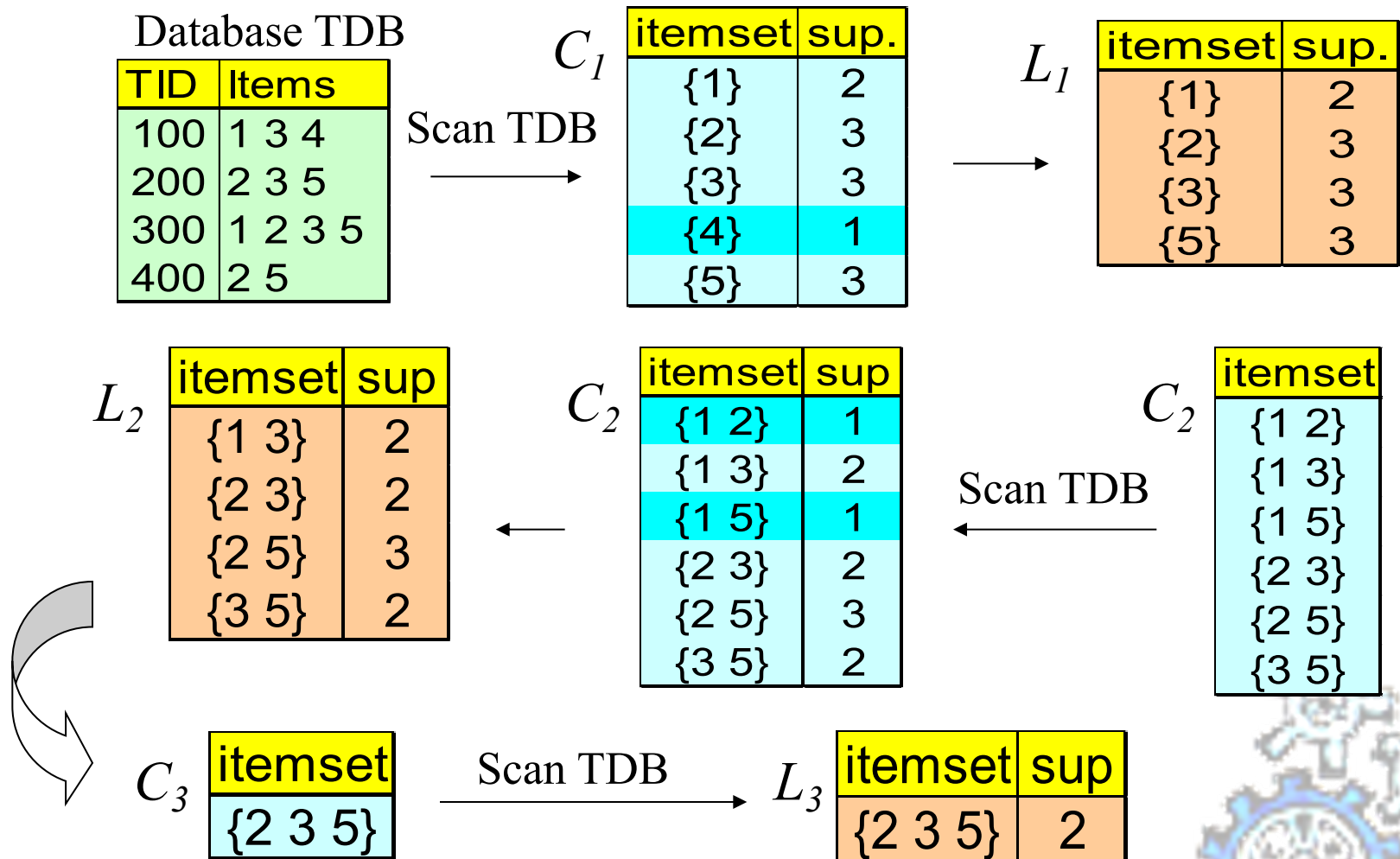


Triplets (3-itemsets)

Itemset	Count
{Bread,Milk,Diaper}	3



# Apriori Execution Example ( $min\_sup = 2$ )



# The Apriori Algorithm

- **Join Step:**  $C_k$  is generated by joining  $L_{k-1}$  with itself
- **Prune Step:** Any  $(k-1)$ -itemset that is not frequent cannot be a subset of a frequent  $k$ -itemset
- **Pseudo-code:**

$C_k$ : Candidate itemset of size  $k$   
 $L_k$ : frequent itemset of size  $k$

$L_1 = \{\text{frequent items}\};$

**for** ( $k = 1; L_k \neq \emptyset; k++$ ) **do begin**

$C_{k+1}$  = candidates generated from  $L_k$ ;

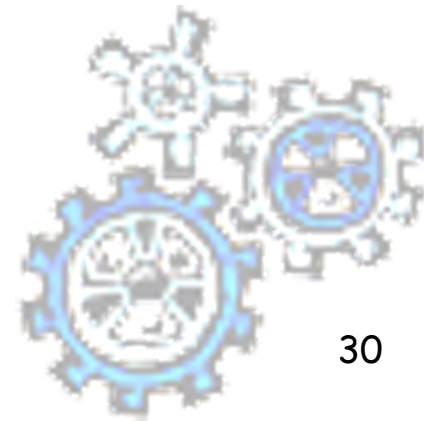
**for each** transaction  $t$  in database **do**

        increment the count of all candidates in  $C_{k+1}$   
        that are contained in  $t$

$L_{k+1}$  = candidates in  $C_{k+1}$  with min\_support

**end**

**return**  $\cup_k L_k$ ;



# How to Generate Candidates?

- Suppose the items in  $L_{k-1}$  are listed in an order
- Step 1: self-joining  $L_{k-1}$ 
  - insert into  $C_k$
  - select  $p.item_1, p.item_2, \dots, p.item_{k-1}, q.item_{k-1}$
  - from  $L_{k-1} p, L_{k-1} q$
  - where  $p.item_1 = q.item_1, \dots, p.item_{k-2} = q.item_{k-2}, p.item_{k-1} < q.item_{k-1}$
- Step 2: pruning
  - forall *itemsets*  $c$  in  $C_k$  do
  - forall  $(k-1)$ -subsets  $s$  of  $c$  do
  - if ( $s$  is not in  $L_{k-1}$ ) then delete  $c$  from  $C_k$



# Example of Generating Candidates

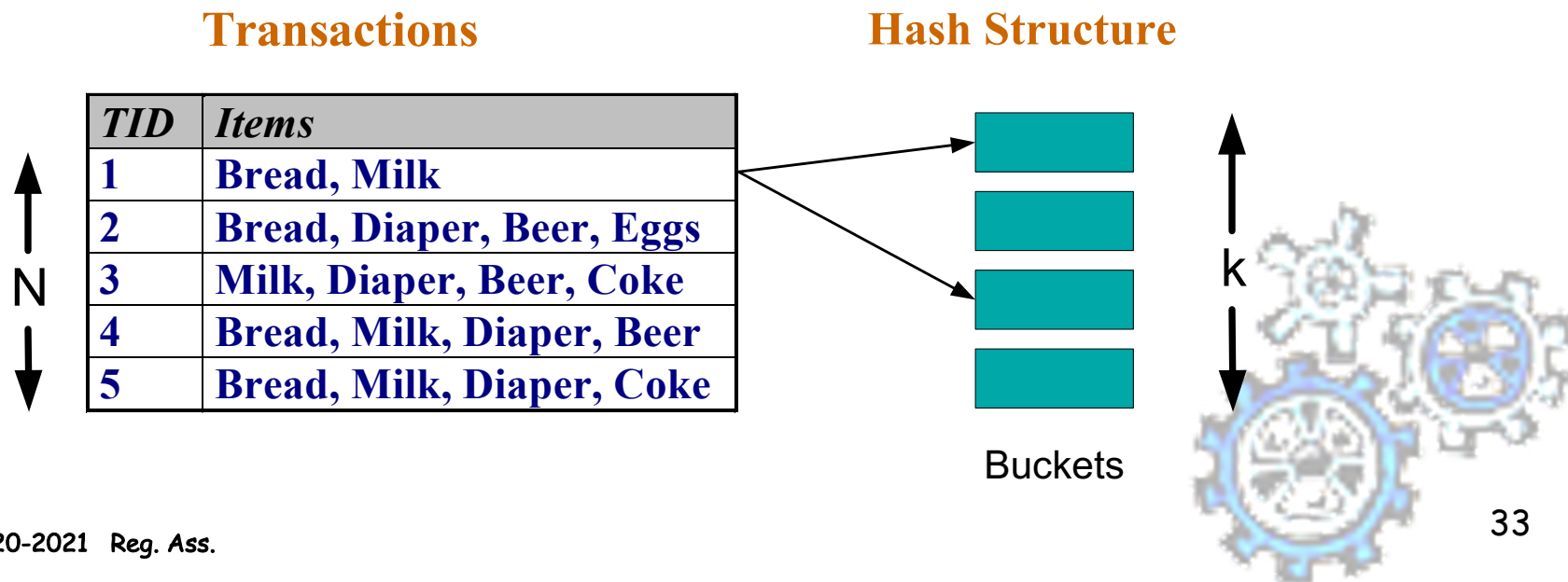
- $L_3 = \{abc, abd, acd, ace, bcd\}$
- Self-joining:  $L_3 * L_3$ 
  - $abcd$  from  $abc$  and  $abd$
  - $acde$  from  $acd$  and  $ace$
- Pruning:
  - $acde$  is removed because  $ade$  is not in  $L_3$
- $C_4 = \{abcd\}$



# Reducing Number of Comparisons

## ■ Candidate counting:

- Scan the database of transactions to determine the support of each candidate itemset
- To reduce the number of comparisons, store the candidates in a hash structure
  - ✓ Instead of matching each transaction against every candidate, match it against candidates contained in the hashed buckets



# Optimizations

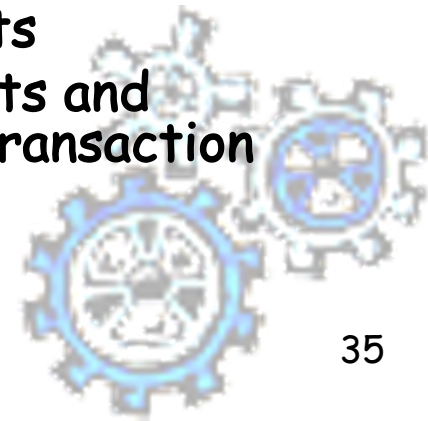
- DHP: Direct Hash and Pruning (Park, Chen and Yu, SIGMOD'95).
- Partitioning Algorithm (Savasere, Omiecinski and Navathe, VLDB'95).
- Sampling (Toivonen'96).
- Dynamic Itemset Counting (Brin et. al. SIGMOD'97)



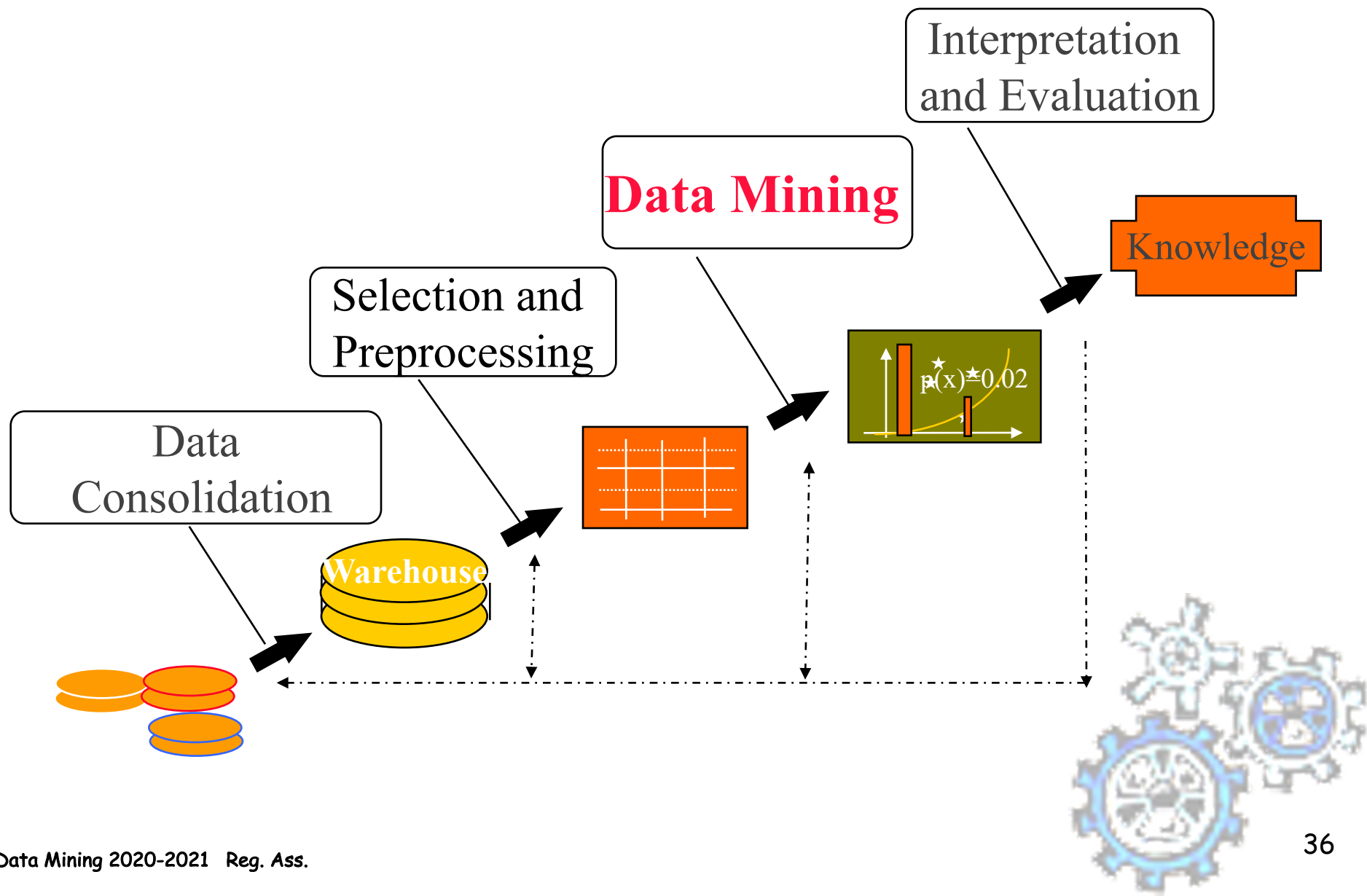


# Factors Affecting Complexity

- **Choice of minimum support threshold**
  - lowering support threshold results in more frequent itemsets
  - this may increase number of candidates and max length of frequent itemsets
- **Dimensionality (number of items) of the data set**
  - more space is needed to store support count of each item
  - if number of frequent items also increases, both computation and I/O costs may also increase
- **Size of database**
  - since Apriori makes multiple passes, run time of algorithm may increase with number of transactions
- **Average transaction width**
  - transaction width increases with denser data sets
  - This may increase max length of frequent itemsets and traversals of hash tree (number of subsets in a transaction increases with its width)



# The KDD process



# Generating Association Rules from Frequent Itemsets

- Only strong association rules are generated
- Frequent itemsets satisfy minimum support threshold
- Strong rules are those that satisfy minimum confidence threshold

$$\frac{\text{support}(A \cup B)}{\text{support}(A)}$$

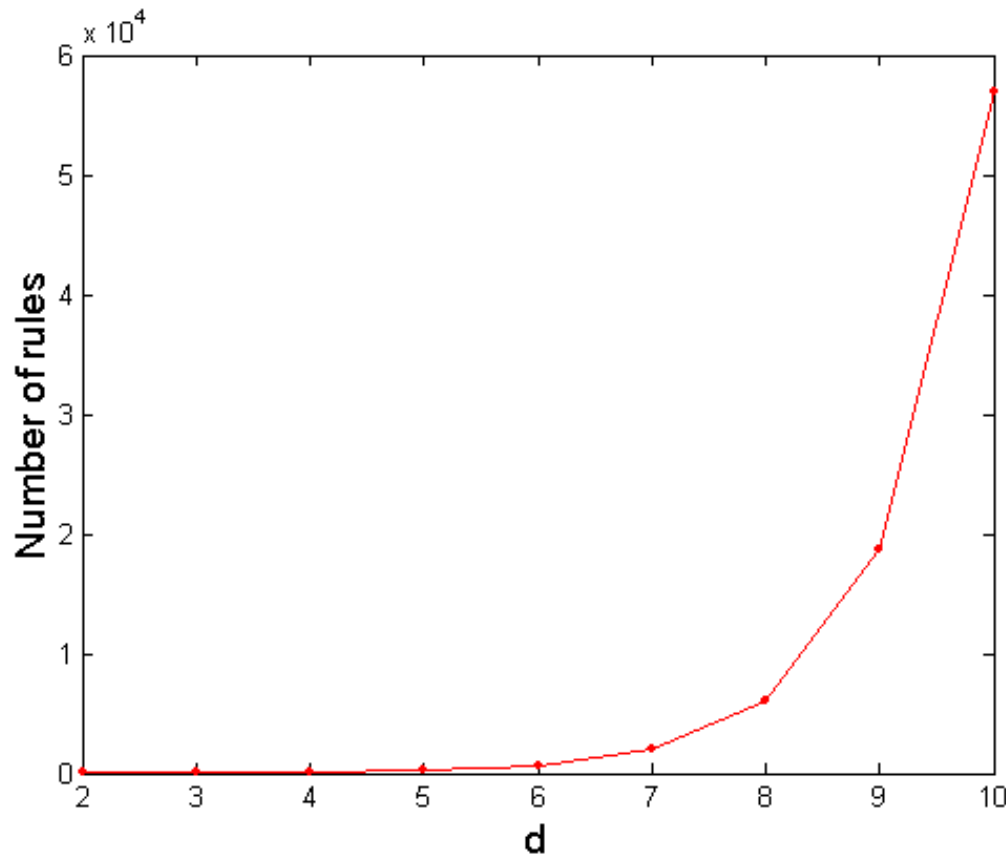
- $\text{confidence}(A \Rightarrow B) = \Pr(B | A) =$

**For each** frequent itemset, **f**, generate all non-empty subsets of **f**  
**For every** non-empty subset **s** of **f** **do**  
    **if**  $\text{support}(\mathbf{f}) / \text{support}(\mathbf{s}) \geq \text{min\_confidence}$  **then**  
        output rule  $\mathbf{s} \Rightarrow (\mathbf{f} - \mathbf{s})$   
**end**

# Computational Complexity

## ■ Given $d$ unique items:

- Total number of itemsets =  $2^d$
- Total number of possible association rules:



$$R = \sum_{k=1}^{d-1} \left[ \binom{d}{k} \times \sum_{j=1}^{d-k} \binom{d-k}{j} \right]$$
$$= 3^d - 2^{d+1} + 1$$

If  $d=6$ ,  $R = 602$  rules



# Rule Generation

- Given a frequent itemset  $L$ , find all non-empty subsets  $f \subset L$  such that  $f \rightarrow L - f$  satisfies the minimum confidence requirement

- If  $\{A,B,C,D\}$  is a frequent itemset, candidate rules:

$ABC \rightarrow D,$	$ABD \rightarrow C,$	$ACD \rightarrow B,$	$BCD \rightarrow A,$
$A \rightarrow BCD,$	$B \rightarrow ACD,$	$C \rightarrow ABD,$	$D \rightarrow ABC$
$AB \rightarrow CD,$	$AC \rightarrow BD,$	$AD \rightarrow BC,$	$BC \rightarrow AD,$
$BD \rightarrow AC,$	$CD \rightarrow AB,$		

- If  $|L| = k$ , then there are  $2^k - 2$  candidate association rules (ignoring  $L \rightarrow \emptyset$  and  $\emptyset \rightarrow L$ )



# Rule Generation

## ■ How to efficiently generate rules from frequent itemsets?

- In general, confidence does not have an anti-monotone property

$c(ABC \rightarrow D)$  can be larger or smaller than  $c(AB \rightarrow D)$

- But confidence of rules generated from the same itemset has an anti-monotone property
- e.g.,  $L = \{A, B, C, D\}$ :

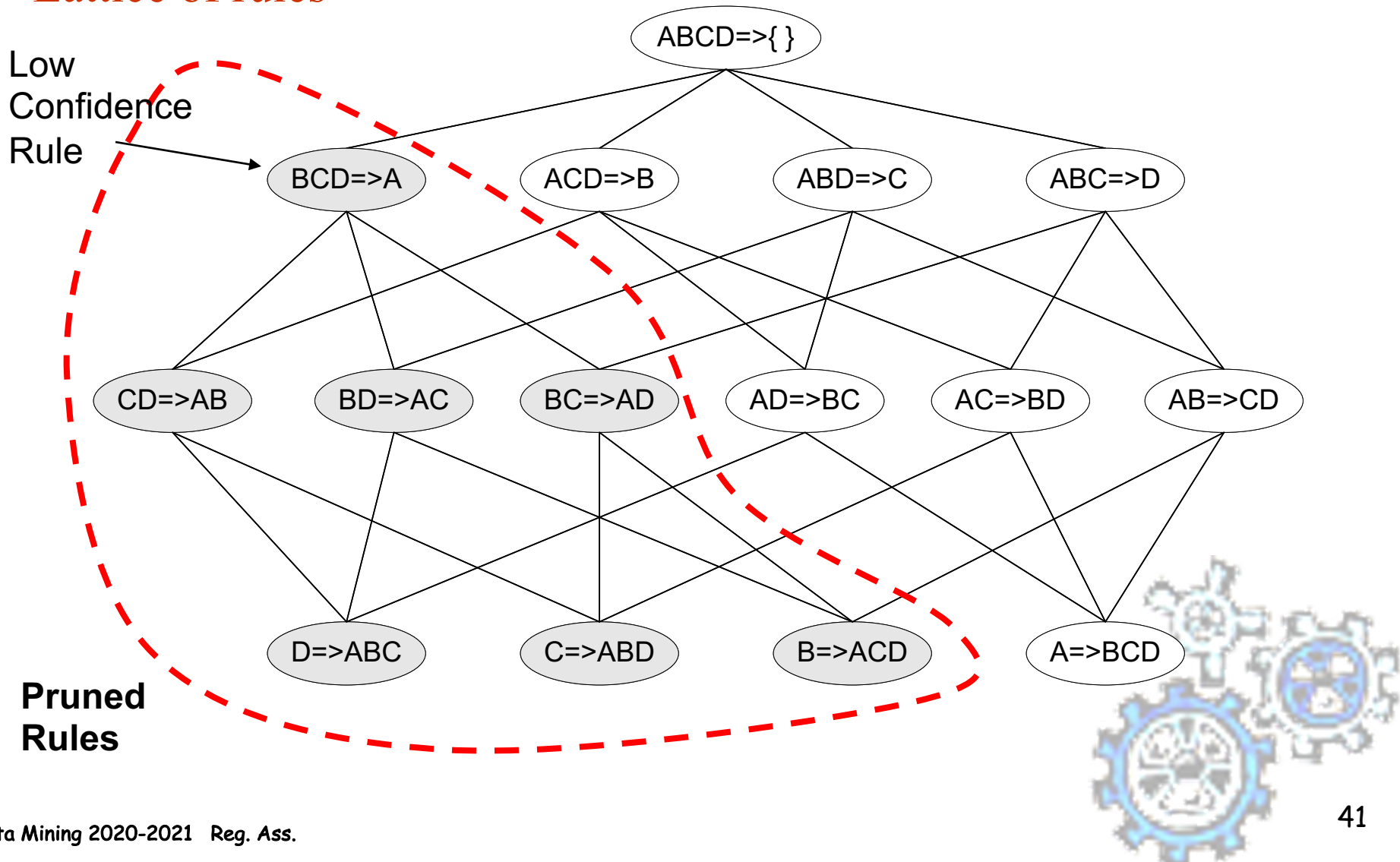
$$c(ABC \rightarrow D) \geq c(AB \rightarrow CD) \geq c(A \rightarrow BCD)$$

- ✓ Confidence is anti-monotone w.r.t. number of items on the RHS of the rule



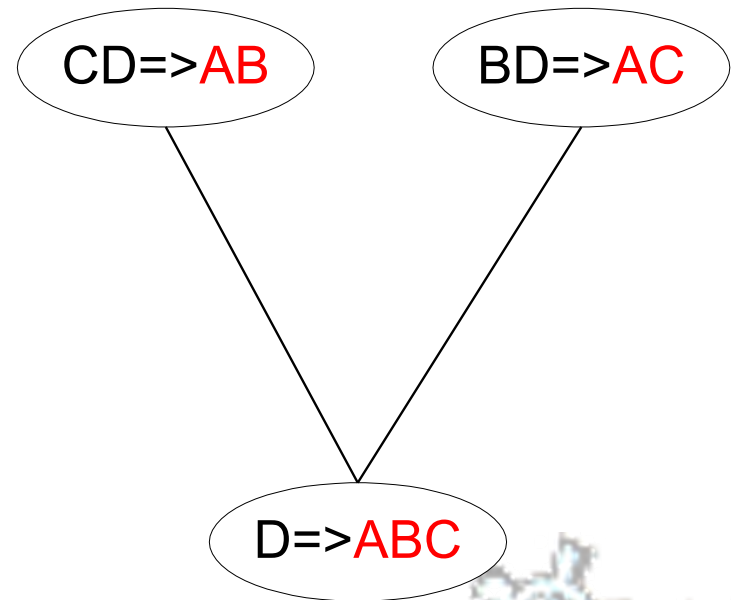
# Rule Generation for Apriori Algorithm

## Lattice of rules



# Rule Generation for Apriori Algorithm

- Candidate rule is generated by merging two rules that share the same prefix in the rule consequent
- $\text{join}(CD \Rightarrow AB, BD \Rightarrow AC)$  would produce the candidate rule  $D \Rightarrow ABC$
- Prune rule  $D \Rightarrow ABC$  if its subset  $AD \Rightarrow BC$  does not have high confidence





# Beyond Support and Confidence

## ■ Example 1: (Aggarwal & Yu, PODS98)

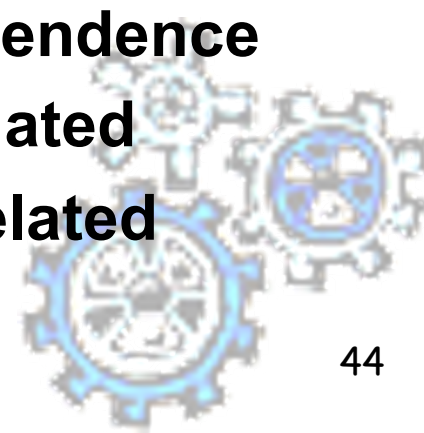
	coffee	not coffee	sum(row)
tea	20	5	25
not tea	70	5	75
sum(col.)	90	10	100

- $\{\text{tea}\} \Rightarrow \{\text{coffee}\}$  has high support (20%) and confidence (80%)
- However, a priori probability that a customer buys coffee is 90%
  - A customer who is known to buy tea is less likely to buy coffee (by 10%)
  - There is a negative correlation between buying tea and buying coffee
  - $\{\sim\text{tea}\} \Rightarrow \{\text{coffee}\}$  has higher confidence(93%)



# Statistical Independence

- Population of 1000 students
  - 600 students know how to swim (S)
  - 700 students know how to bike (B)
  - 420 students know how to swim and bike (S,B)
- $P(S \wedge B) = 420/1000 = 0.42$
- $P(S) \times P(B) = 0.6 \times 0.7 = 0.42$
- $P(S \wedge B) = P(S) \times P(B) \Rightarrow$  Statistical independence
- $P(S \wedge B) > P(S) \times P(B) \Rightarrow$  Positively correlated
- $P(S \wedge B) < P(S) \times P(B) \Rightarrow$  Negatively correlated



# Correlation and Interest

- Two events are independent if  $P(A \wedge B) = P(A) * P(B)$ , otherwise are correlated.
- Interest =  $P(A \wedge B) / P(B) * P(A)$
- Interest expresses measure of correlation
  - $= 1 \Rightarrow A$  and  $B$  are independent events
  - **less than 1**  $\Rightarrow A$  and  $B$  negatively correlated,
  - **greater than 1**  $\Rightarrow A$  and  $B$  positively correlated.
  - In our example,  $I(\text{buy tea} \wedge \text{buy coffee}) = 0.89$  i.e. they are negatively correlated.



# Computing Interestingness Measure

- Given a rule  $X \rightarrow Y$ , information needed to compute rule interestingness can be obtained from a contingency table

Contingency table for  $X \rightarrow Y$

	$Y$	$\overline{Y}$	
$X$	$f_{11}$	$f_{10}$	$f_{1+}$
$\overline{X}$	$f_{01}$	$f_{00}$	$f_{0+}$
	$f_{+1}$	$f_{+0}$	$ T $

$f_{11}$ : support of  $X$  and  $Y$

$f_{10}$ : support of  $X$  and  $\overline{Y}$

$f_{01}$ : support of  $\overline{X}$  and  $Y$

$f_{00}$ : support of  $\overline{X}$  and  $\overline{Y}$

Used to define various measures

- ◆ support, confidence, lift, Gini, J-measure, etc.

# ■ Wrap up



# Frequent Itemsets

Transaction ID	Items Bought
1	dairy,fruit
2	dairy,fruit, vegetable
3	dairy
4	fruit, cereals

Support({dairy}) = 3/4 (75%)

Support({fruit}) = 3/4 (75%)

Support({dairy, fruit}) = 2/4 (50%)

If minsup = 60%, then

{dairy} and {fruit} are frequent while {dairy, fruit} is not.



# Association Rules: Measures

- Let  $A$  and  $B$  be a partition of an itemset  $I$  :

$$A \Rightarrow B [s, c]$$

$A$  and  $B$  are itemsets

$$s = \text{support of } A \Rightarrow B = \text{support}(A, B)$$

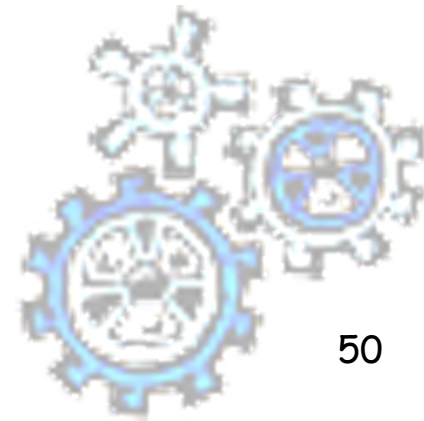
$$c = \text{confidence of } A \Rightarrow B = \text{support}(A, B) / \text{support}(A)$$

- Measure for rules:

- ✓ minimum support  $\sigma$

- ✓ minimum confidence  $\gamma$

- The rules holds if :  $s \geq \sigma$  and  $c \geq \gamma$



# Association Rules: Meaning

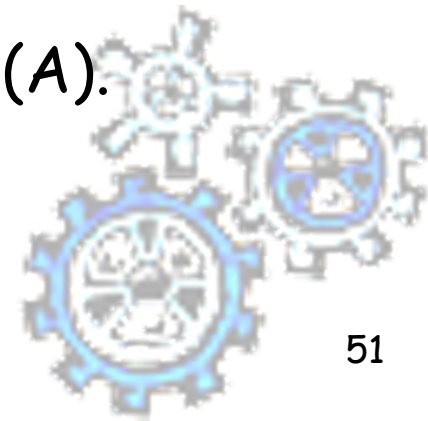
$$A \Rightarrow B [s, c]$$

**Support:** denotes the frequency of the rule within transactions. A high value means that the rule involve a great part of database.

$$\text{support}(A \Rightarrow B) = p(A \ \& \ B)$$

**Confidence:** denotes the percentage of transactions containing A which contain also B. It is an estimation of conditioned probability .

$$\text{confidence}(A \Rightarrow B) = p(B|A) = p(A \ \& \ B)/p(A).$$





# Association Rules - the parameters $\sigma$ and $\gamma$

Minimum Support  $\sigma$  :

High  $\Rightarrow$  **few** frequent itemsets  
 $\Rightarrow$  **few** valid rules which occur **very often**

Low  $\Rightarrow$  **many** valid rules which occur **rarely**

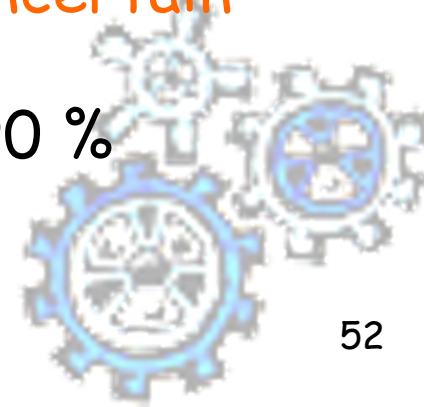
Minimum Confidence  $\gamma$  :

High  $\Rightarrow$  **few** rules, but all “**almost logically true**”

Low  $\Rightarrow$  many rules, but many of them very “**uncertain**”

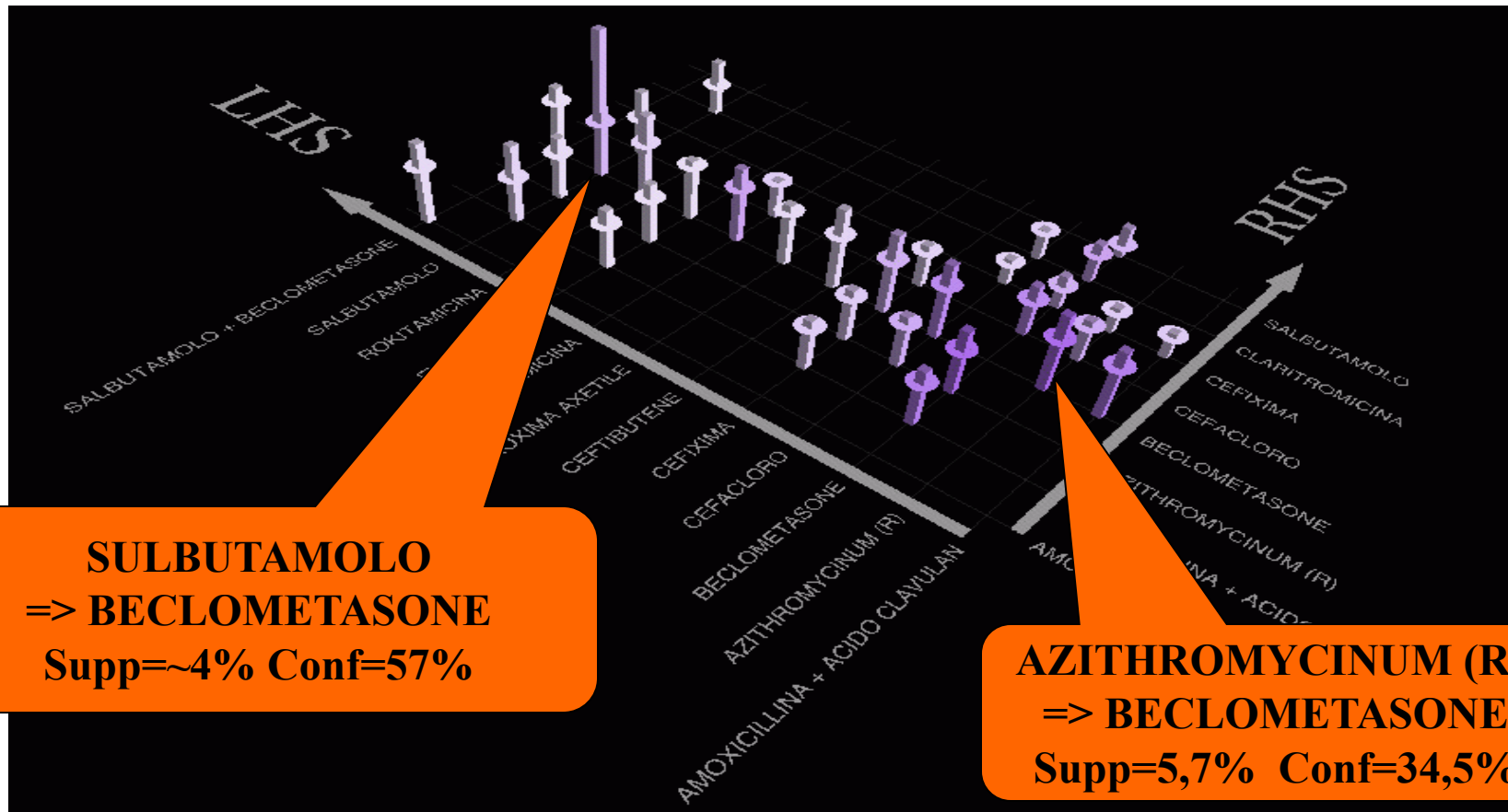
Typical Values:  $\sigma = 2 \div 10 \%$

$\gamma = 70 \div 90 \%$



# Association Rules - visualization

(Patients <15 old for USL 19 (a unit of Sanitary service),  
January-September 1997)



# Association Rules - bank transactions

Step 1: Create groups of customers (cluster) on the base of demographical data.

Step 2: Describe customers of each cluster by mining association rules.

Example:

Rules on cluster 6  
(23,7% of dataset):

Group	Support	Confidence	Body	Head
1	0.277	91.4	1.3	[TERM DEPOSITS] AND [BUSINESS CREDIT CARD] AND [TELEBANKING] AND [BUSINESS SAVINGS]
1	0.164	86.4	1.3	[SAVINGS] AND [ATH CARD] AND [BUSINESS CREDIT CARD] AND [TELEBANKING] AND [BUSINESS SAVINGS]
1	0.104	85.7	1.9	[SAVINGS] AND [INTERNET BANKING] AND [LEASES]
1	0.138	84.2	1.2	[PERSONAL BANKING] AND [TERM DEPOSITS] AND [BUSINESS CREDIT CARD] AND [BUSINESS SAVINGS]
1	0.251	82.9	1.2	[SAVINGS] AND [TERM DEPOSITS] AND [ATH CARD] AND [TELEBANKING] AND [BUSINESS SAVINGS]
1	0.328	82.6	1.2	[SAVINGS] AND [ATH CARD] AND [BUSINESS CREDIT CARD] AND [TELEBANKING] AND [BUSINESS SAVINGS]
1	0.242	82.4	1.2	[PERSONAL BANKING] AND [TERM DEPOSITS] AND [BUSINESS SAVINGS]
1	0.631	81.1	1.2	[SAVINGS] AND [BUSINESS CREDIT CARD] AND [TELEBANKING] AND [BUSINESS SAVINGS]
1	0.138	80.8	1.2	[SAVINGS] AND [ATH CARD] AND [BUSINESS CREDIT CARD] AND [TELEBANKING] AND [INTERNET BANKING] AND [BUSINESS SAVINGS]
1	0.138	80.8	1.2	[SAVINGS] AND [TERM DEPOSITS] AND [TEL
1	0.458	79.1	1.2	[SAVINGS] AND [TERM DEPOSITS] AND [TELEBANKING] AND [BUSINESS SAVINGS]
1	0.130	78.9	1.2	[SAVINGS] AND [PERSONAL BANKING] AND [BUSINESS CREDIT CARD] AND [TELEBANKING] AND [BUSINESS SAVINGS]
1	0.346	78.4	1.2	[SAVINGS] AND [PERSONAL BANKING] AND [BUSINESS CREDIT CARD] AND [BUSINESS SAVINGS]
1	1.037	77.9	1.1	[SAVINGS] AND [TERM DEPOSITS] AND [ATH CARD] AND [BUSINESS CREDIT CARD] AND [TELEBANKING] AND [INTERNET BANKING]
1	0.182	77.8	1.7	[SAVINGS] AND [TERM DEPOSITS] AND [ATH CARD] AND [INTERNET BANKING] AND [BUSINESS SAVINGS]
				[BUSINESS CREDIT CARD]

# Cluster 6 (23.7% of customers)

Group	Support	Confidence	Body	
1	0.277	91.4	1.3	==> Head [TERM DEPOSITS] AND [BUSINESS CREDIT CARD] AND [TELEBANKING] AND [BUSINESS SAVINGS] ==> [SAVINGS]
1	0.164	86.4	1.3	[TERM DEPOSITS] AND [ATH CARD] AND [BUSINESS CREDIT CARD] AND [TELEBANKING] AND [BUSINESS SAVINGS] ==> [SAVINGS]
1	0.104	85.7	1.9	[SAVINGS] AND [INTERNET BANKING] AND [LEASES] ==> [TELEBANKING]
1	0.138	84.2	1.2	[PERSONAL BANKING] AND [TERM DEPOSITS] AND [BUSINESS CREDIT CARD] AND [BUSINESS SAVINGS] ==> [SAVINGS]
1	0.251	82.9	1.2	[TERM DEPOSITS] AND [ATH CARD] AND [TELEBANKING] AND [BUSINESS SAVINGS] ==> [SAVINGS]
1	0.328	82.6	1.2	[ATH CARD] AND [BUSINESS CREDIT CARD] AND [TELEBANKING] AND [BUSINESS SAVINGS] ==> [SAVINGS]
1	0.242	82.4	1.2	[PERSONAL BANKING] AND [TERM DEPOSITS] AND [BUSINESS SAVINGS] ==> [SAVINGS]
1	0.631	81.1	1.2	[BUSINESS CREDIT CARD] AND [TELEBANKING] AND [BUSINESS SAVINGS] ==> [SAVINGS]
1	0.138	80.8	1.2	[ATH CARD] AND [BUSINESS CREDIT CARD] AND [TELEBANKING] AND [INTERNET BANKING] AND [BUSINESS SAVINGS] ==> [SAVINGS]
1	0.138	80.8	1.2	[TERM DEPOSITS] AND [TEL --> [SAVINGS]
1	0.458	79.1	1.2	[TERM DEPOSITS] AND [TELEBANKING] AND [BUSINESS SAVINGS] ==> [SAVINGS]
1	0.130	78.9	1.2	[PERSONAL BANKING] AND [BUSINESS CREDIT CARD] AND [TELEBANKING] AND [BUSINESS SAVINGS] ==> [SAVINGS]
1	0.346	78.4	1.2	[PERSONAL BANKING] AND [BUSINESS CREDIT CARD] AND [BUSINESS SAVINGS] ==> [SAVINGS]
1	1.037	77.9	1.1	[TERM DEPOSITS] AND [ATH CARD] AND [BUSINESS CREDIT CARD] AND [TELEBANKING] AND [INTERNET BANKING] ==> [SAVINGS]
1	0.182	77.8	1.7	[TERM DEPOSITS] AND [ATH CARD] AND [INTERNET BANKING] AND [BUSINESS SAVINGS] --> [BUSINESS CREDIT CARD]

# Table (6.1)

Table 6.1. Example of market basket transactions.

Customer ID	Transaction ID	Items Bought
1	0001	$\{a, d, e\}$
1	0024	$\{a, b, c, e\}$
2	0012	$\{a, b, d, e\}$
2	0031	$\{a, c, d, e\}$
3	0015	$\{b, c, e\}$
3	0022	$\{b, d, e\}$
4	0029	$\{c, d\}$
4	0040	$\{a, b, c\}$
5	0033	$\{a, d, e\}$
5	0038	$\{a, b, e\}$

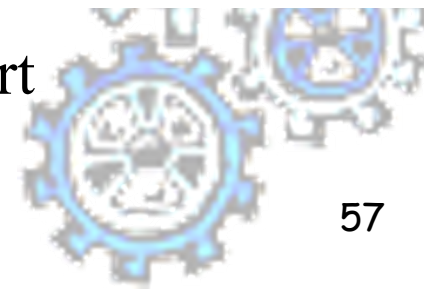
Support?: e, (b,d), (b,d,e)



**Table 6.2. Market basket transactions.**

Transaction ID	Items Bought
1	{Milk, Beer, Diapers}
2	{Bread, Butter, Milk}
3	{Milk, Diapers, Cookies}
4	{Bread, Butter, Cookies}
5	{Beer, Cookies, Diapers}
6	{Milk, Diapers, Bread, Butter}
7	{Bread, Butter, Diapers}
8	{Beer, Diapers}
9	{Milk, Diapers, Bread, Butter}
10	{Beer, Cookies}

Max size of itemset, 2-itemsets with larger support



Most FIs  
 Fewest FIs  
 Longest  
 Highest Maximum support

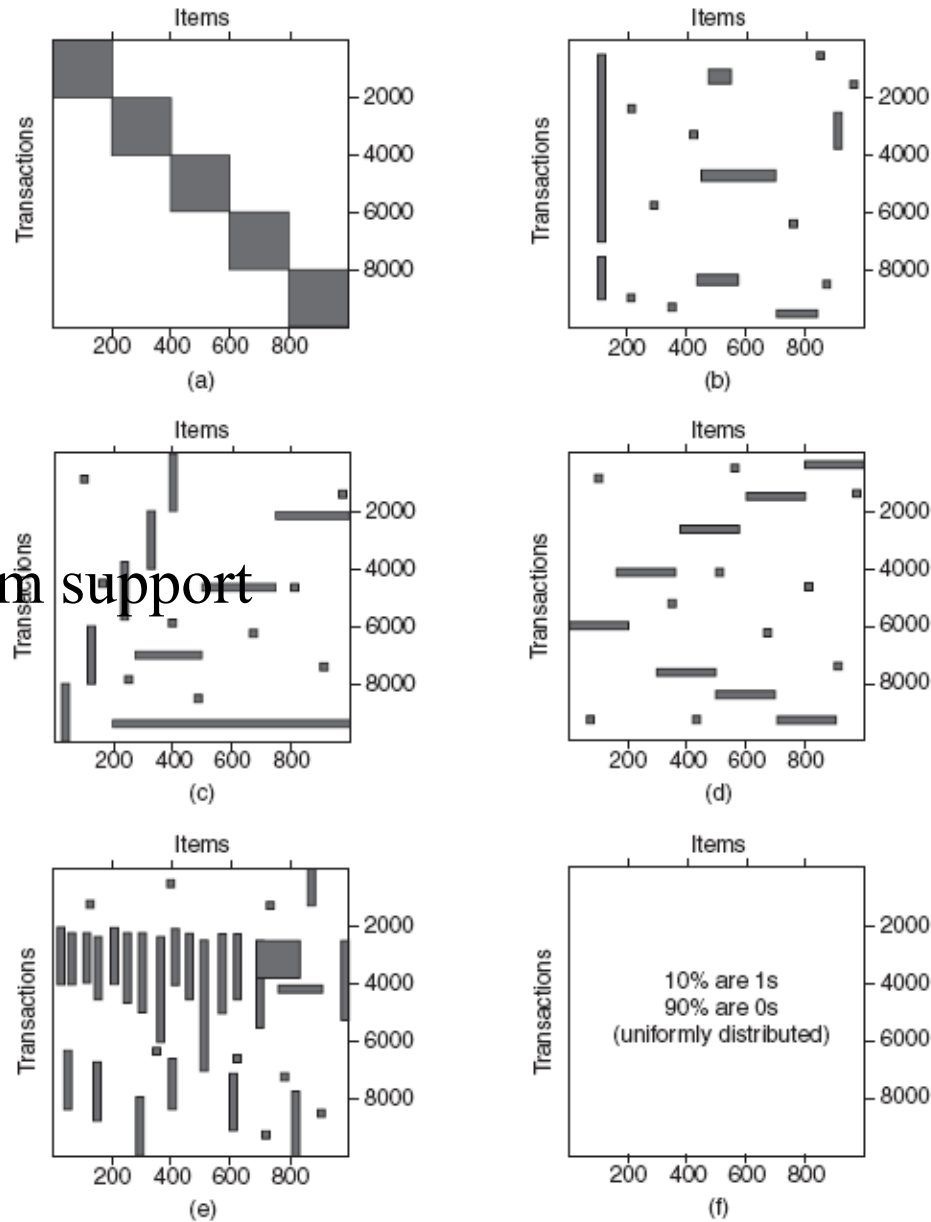


Figure 6.6. Figures for Exercise 14.



# Factors Affecting Complexity

---

- Choice of minimum support threshold
  - lowering support threshold results in more frequent itemsets
  - this may increase number of candidates and max length of frequent itemsets
- Dimensionality (number of items) of the data set
  - more space is needed to store support count of each item
  - if number of frequent items also increases, both computation and I/O costs may also increase
- Size of database
  - since Apriori makes multiple passes, run time of algorithm may increase with number of transactions
- Average transaction width
  - transaction width increases with denser data sets
  - This may increase max length of frequent itemsets and traversals of hash tree (number of subsets in a transaction increases with its width)



# Compact Representation of Frequent Itemsets

- Some itemsets are redundant because they have identical support as their supersets

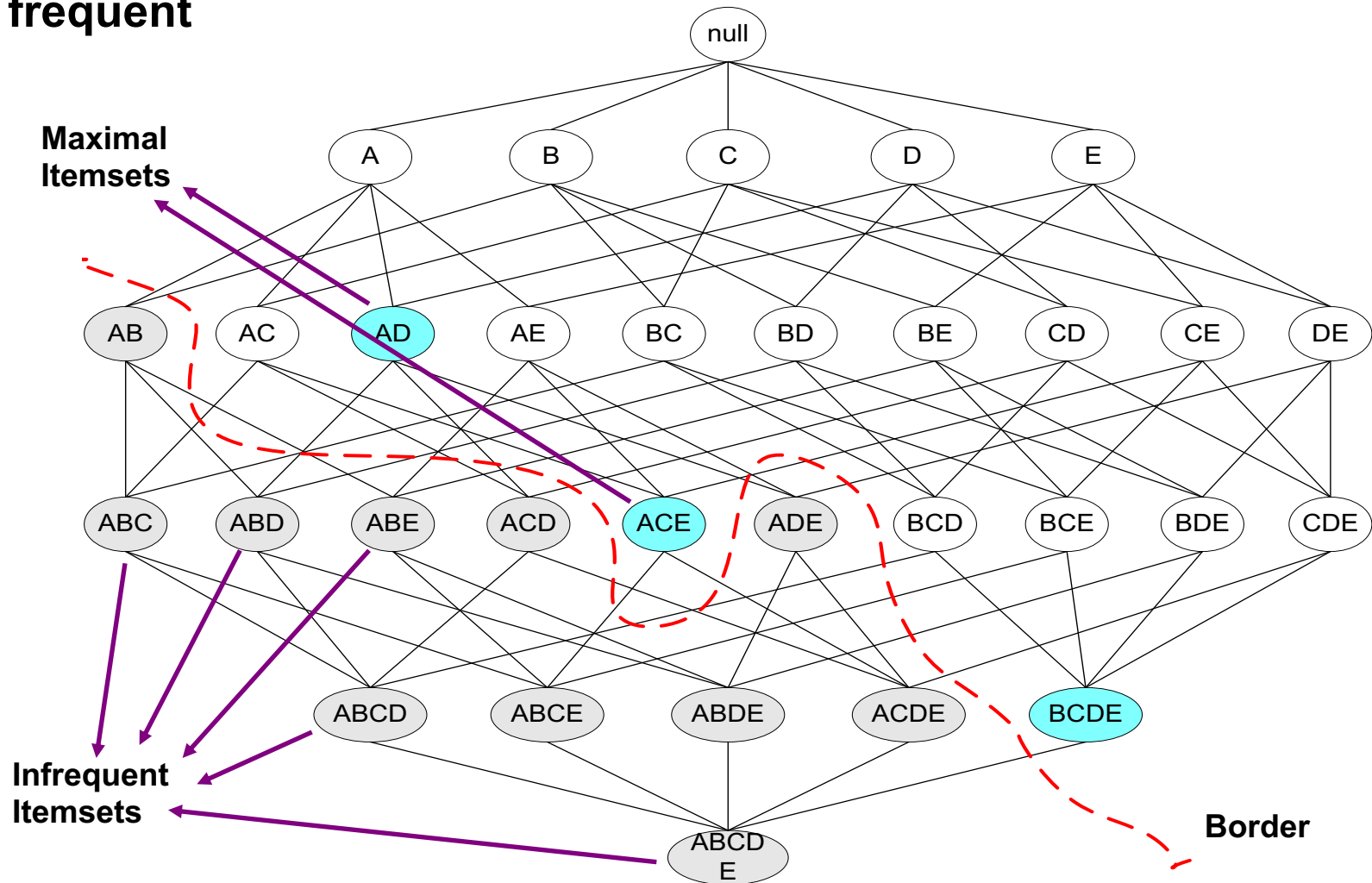
TID	A1	A2	A3	A4	A5	A6	A7	A8	A9	A10	B1	B2	B3	B4	B5	B6	B7	B8	B9	B10	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10
1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1

- Number of frequent itemsets =  $3 \times \sum_{k=1}^{10} \binom{10}{k}$

- Need a compact representation

# Maximal Frequent Itemset

An itemset is maximal frequent if none of its immediate supersets is frequent



# Closed Itemset

- An itemset is closed if none of its immediate supersets has the same support as the itemset

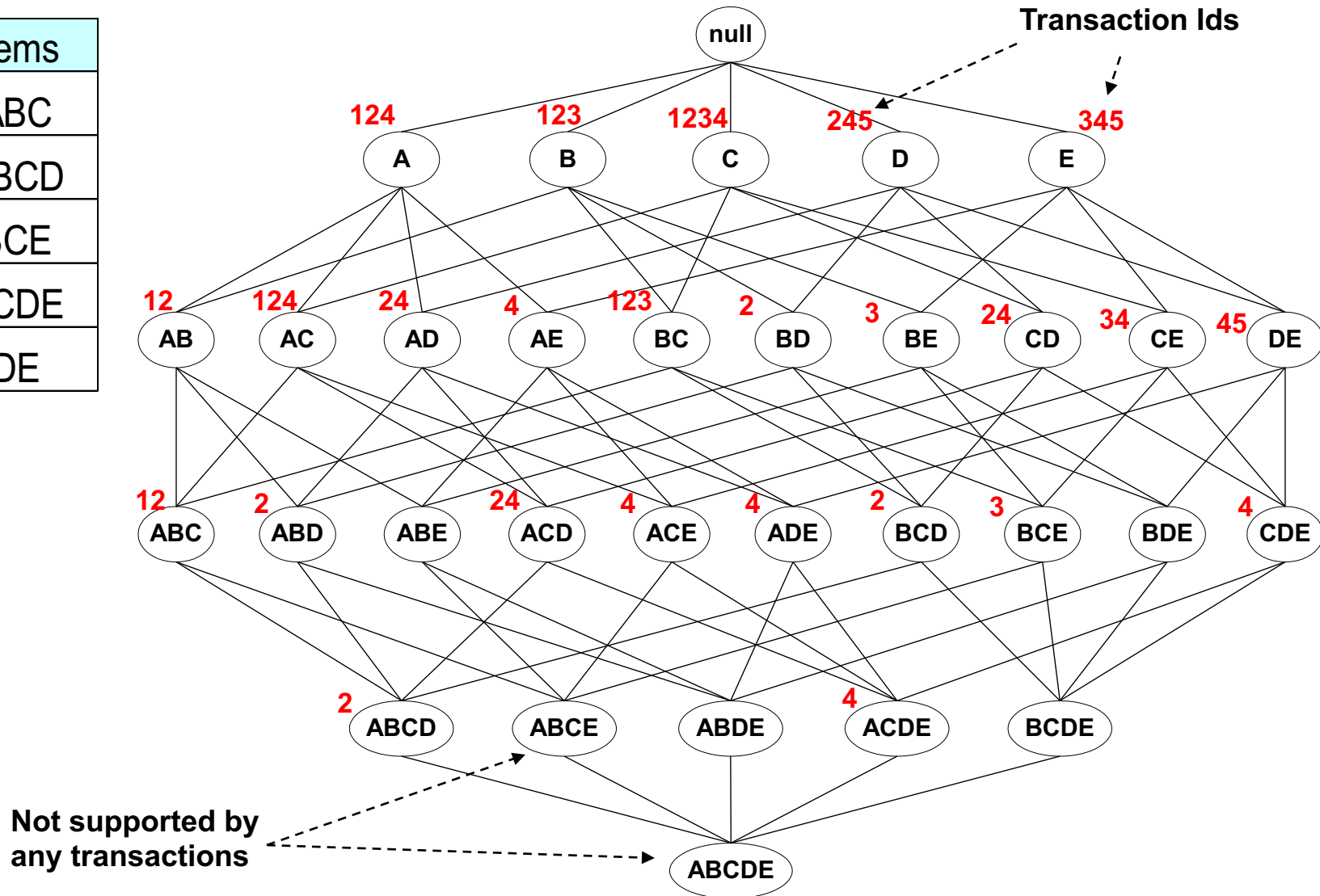
TID	Items
1	{A,B}
2	{B,C,D}
3	{A,B,C,D}
4	{A,B,D}
5	{A,B,C,D}

Itemset	Support
{A}	4
{B}	5
{C}	3
{D}	4
{A,B}	4
{A,C}	2
{A,D}	3
{B,C}	3
{B,D}	4
{C,D}	3

Itemset	Support
{A,B,C}	2
{A,B,D}	3
{A,C,D}	2
{B,C,D}	3
{A,B,C,D}	2

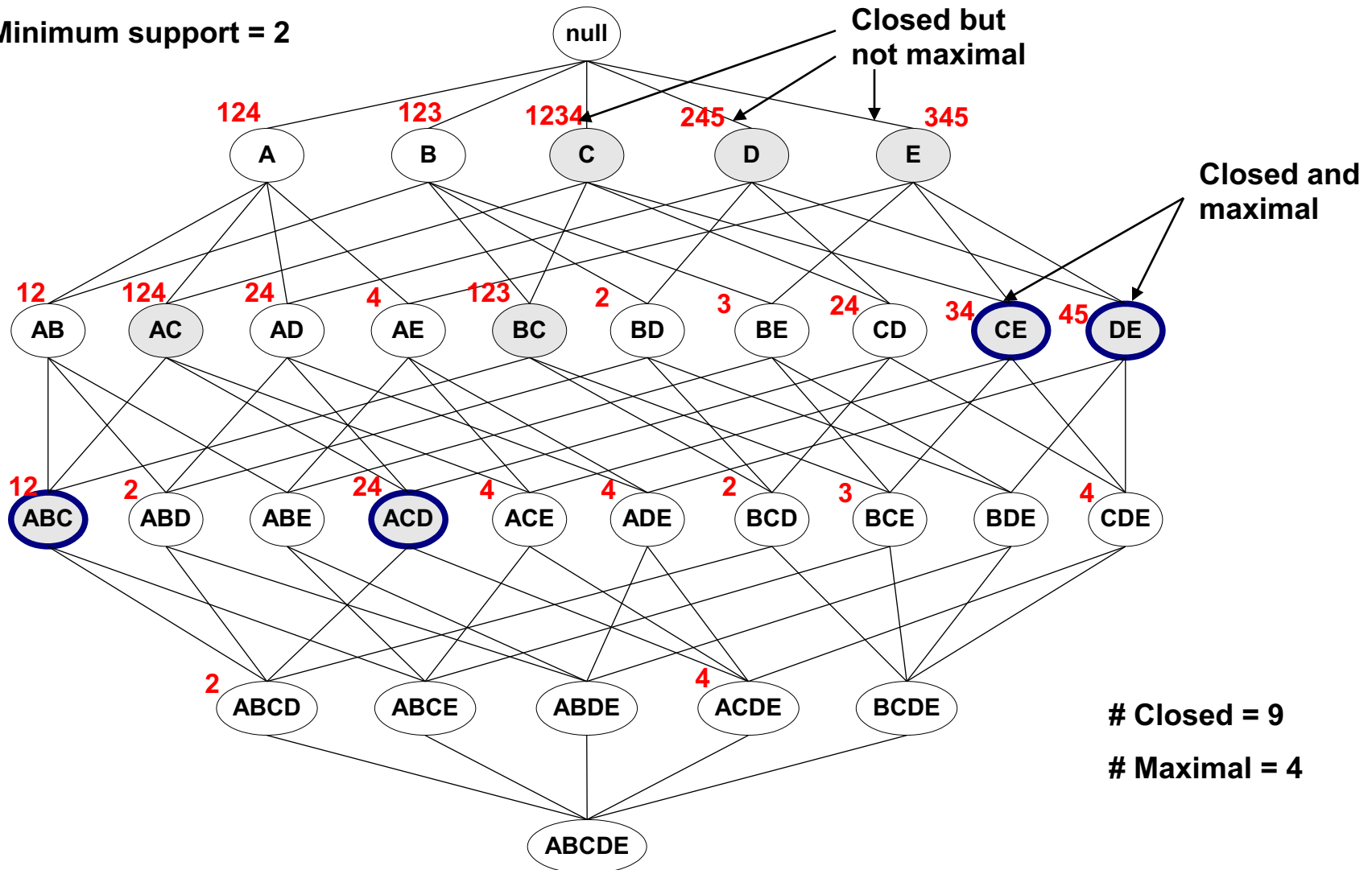
# Maximal vs Closed Itemsets

TID	Items
1	ABC
2	ABCD
3	BCE
4	ACDE
5	DE



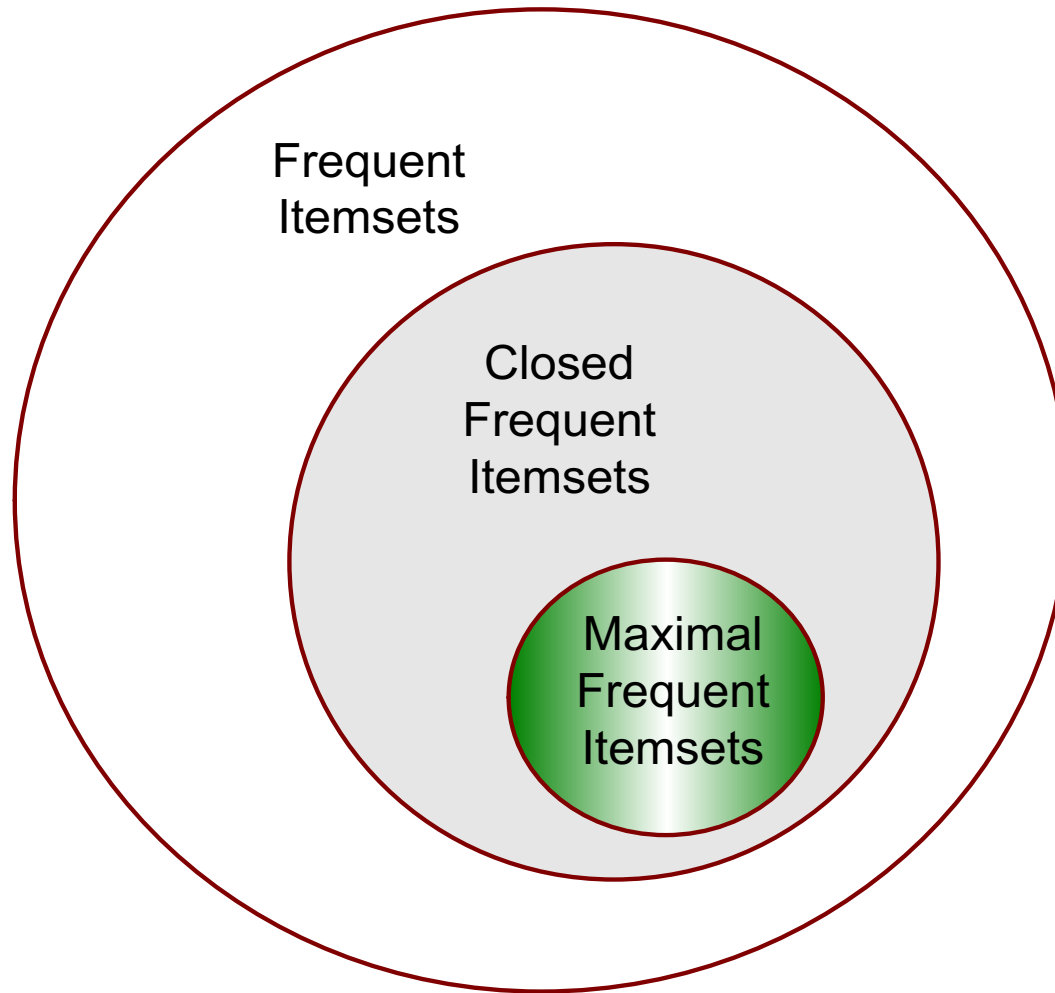
# Maximal vs Closed Frequent Itemsets

Minimum support = 2



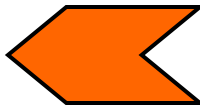
# Maximal vs Closed Itemsets

---



# Association rules - module outline

- What are association rules (AR) and what are they used for:
  - The paradigmatic application: Market Basket Analysis
  - The single dimensional AR (intra-attribute)
- How to compute AR
  - Basic Apriori Algorithm and its optimizations
  - Multi-Dimension AR (inter-attribute)
  - Quantitative AR
  - Constrained AR
- How to reason on AR and how to evaluate their quality
  - Multiple-level AR
  - Interestingness
  - Correlation vs. Association



# Multidimensional AR

Associations between values of different attributes :

CID	nationality	age	income
1	Italian	50	low
2	French	40	high
3	French	30	high
4	Italian	50	medium
5	Italian	45	high
6	French	35	high

RULES:

**nationality = French**  $\Rightarrow$  **income = high** [50%, 100%]

**income = high**  $\Rightarrow$  **nationality = French** [50%, 75%]

**age = 50**  $\Rightarrow$  **nationality = Italian** [33%, 100%]



# Single-dimensional vs multi-dimensional AR

## Single-dimensional (Intra-attribute)

The events are: *items A, B and C belong to the same transaction*

Occurrence of events: *transactions*

## Multi-dimensional (Inter-attribute)

The events are : *attribute A assumes value a, attribute B assumes value b and attribute C assumes value c.*

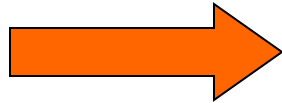
Occurrence of events: *tuples*



# Single-dimensional vs Multi-dimensional AR

## Multi-dimensional

<1, Italian, 50, low>  
<2, French, 45, high>

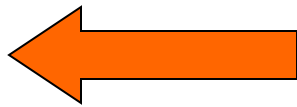


## Single-dimensional

<1, {nat/Ita, age/50, inc/low}>  
<2, {nat/Fre, age/45, inc/high}>

Schema: <ID, a?, b?, c?, d?>

<1, yes, yes, no, no>  
<2, yes, no, yes, no>



<1, {a, b}>  
<2, {a, c}>



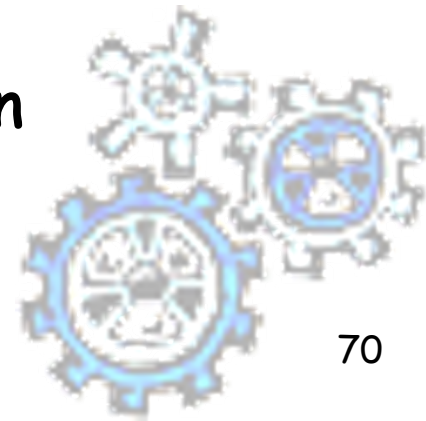
# Quantitative Attributes

- Quantitative attributes (e.g. age, income)
- Categorical attributes (e.g. color of car)

CID	height	weight	income
1	168	75,4	30,5
2	175	80,0	20,3
3	174	70,3	25,8
4	170	65,2	27,0

**Problem:** too many distinct values

**Solution:** transform quantitative attributes in categorical ones via **discretization**.



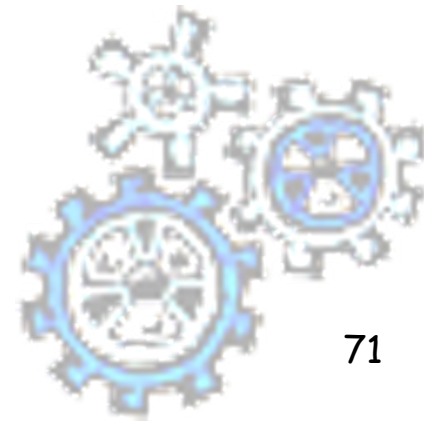
# Quantitative Association Rules

CID	Age	Married	NumCars
1	23	No	1
2	25	Yes	1
3	29	No	0
4	34	Yes	2
5	38	Yes	2

**[Age: 30..39] and [Married: Yes]  $\Rightarrow$  [NumCars:2]**

support = 40%

confidence = 100%



# Discretization of quantitative attributes

**Solution:** each value is replaced by the interval to which it belongs.

height: 0-150cm, 151-170cm, 171-180cm, >180cm

weight: 0-40kg, 41-60kg, 60-80kg, >80kg

income: 0-10ML, 11-20ML, 20-25ML, 25-30ML, >30ML

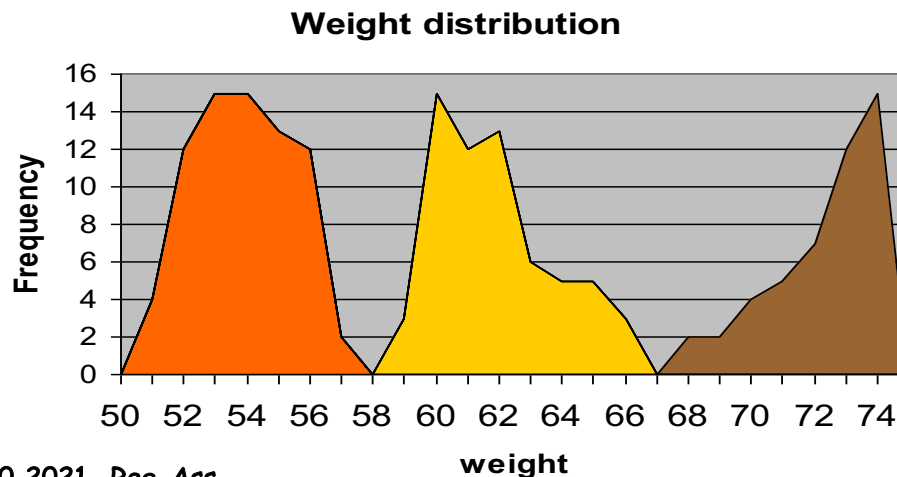
CID	height	weight	income
1	151-171	60-80	>30
2	171-180	60-80	20-25
3	171-180	60-80	25-30
4	151-170	60-80	25-30

**Problem:** the discretization may be useless (see weight).



# How to choose intervals?

1. Interval with a fixed "reasonable" granularity  
Ex. **intervals of 10 cm for height.**
2. Interval size is defined by some domain dependent criterion  
Ex.: **0-20ML, 21-22ML, 23-24ML, 25-26ML, >26ML**
3. Interval size determined by analyzing data, studying the distribution or using clustering



**50 - 58 kg**  
**59-67 kg**  
**> 68 kg**



# Discretization of quantitative attributes

1. Quantitative attributes are **statically** discretized by using predefined concept hierarchies:
  - elementary use of background knowledge

Loose interaction between Apriori and discretizer

2. Quantitative attributes are **dynamically** discretized
  - into “bins” based on the distribution of the data.
  - considering the distance between data points.

Tighter interaction between Apriori and discretizer



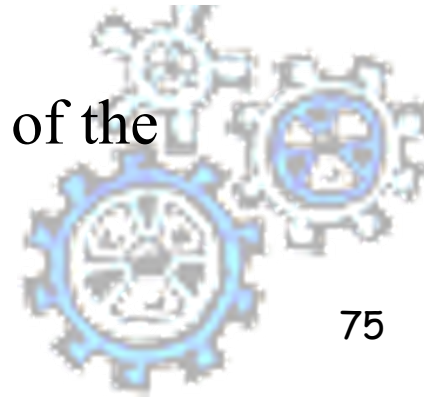
# Quantitative Association Rules

RecordID	Age	Married	NumCars
100	23	No	1
200	25	Yes	1
300	29	No	0
400	34	Yes	2
500	38	Yes	2



Sample Rules	Support	Confidence
<age:30..39> and <married: yes> ==> <numCars:2>	40%	100%
<NumCars: 0..1> ==> <Married: No>	40%	66.70%

Handling quantitative rules may require mapping of the **continuous** variables into **Boolean**





# Mapping Quantitative to Boolean

- One possible solution is to map the problem to the Boolean association rules:

- discretize a non-categorical attribute to intervals, e.g., Age [20,29], [30,39],...
- categorical attributes: each value becomes one item
- non-categorical attributes: each interval becomes one item

- Problems with the mapping

- too few intervals: lost information
- too low support: too many rules

RecordID	Age	Married	NoCars
100	23	No	1
500	38	Yes	2

RecID	Age: 20..29	Age: 30..39	Married: Yes	Married: No	Cars: 0	Cars: 1	Cars: 2
100	1	0	0	1	0	1	0
500	0	1	1	0	0	0	1

# Constraints and AR

- **Preprocessing:** use constraints to focus on a subset of transactions
  - Example: find association rules where the prices of all items are at most 200 Euro
- **Optimizations:** use constraints to optimize Apriori algorithm
  - Anti-monotonicity: when a set violates the constraint, so does any of its supersets.
  - Apriori algorithm uses this property for pruning
- **Push constraints as deep as possible** inside the frequent set computation



# Constraint-based AR

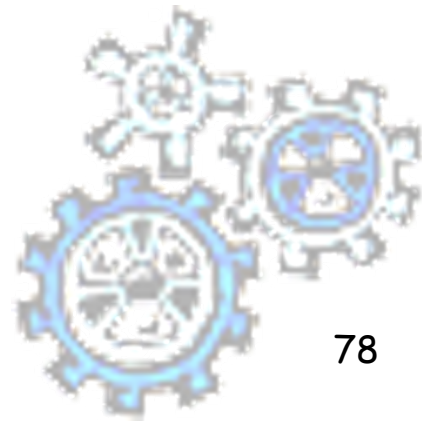
## ■ What kinds of constraints can be used in mining?

### ■ Data constraints:

- ✓ SQL-like queries
  - Find product pairs sold together in **Vancouver** in **Dec.'98**.
- ✓ OLAP-like queries (**Dimension/level**)
  - in relevance to **region, price, brand, customer category**.

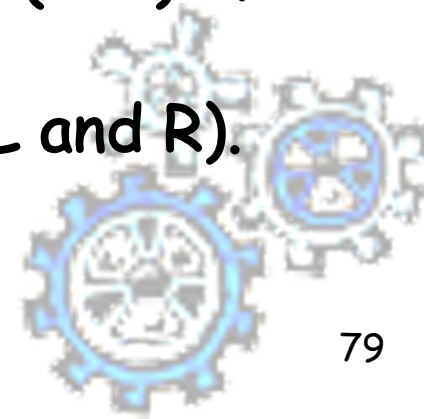
### ■ Rule constraints:

- ✓ specify the form or property of rules to be mined.
- ✓ Constraint-based AR



# Rule Constraints

- Two kind of constraints:
  - Rule form constraints: meta-rule guided mining.
    - ✓  $P(x, y) \wedge Q(x, w) \rightarrow \text{takes}(x, \text{"database systems"})$ .
  - Rule content constraint: constraint-based query optimization (Ng, et al., SIGMOD'98).
    - ✓  $\text{sum}(\text{LHS}) < 100 \wedge \text{min}(\text{LHS}) > 20 \wedge \text{sum}(\text{RHS}) > 1000$
- 1-variable vs. 2-variable constraints (Lakshmanan, et al. SIGMOD'99):
  - 1-var: A constraint confining only one side (L/R) of the rule, e.g., as shown above.
  - 2-var: A constraint confining both sides (L and R).
    - ✓  $\text{sum}(\text{LHS}) < \text{min}(\text{RHS}) \wedge \text{max}(\text{RHS}) < 5 * \text{sum}(\text{LHS})$



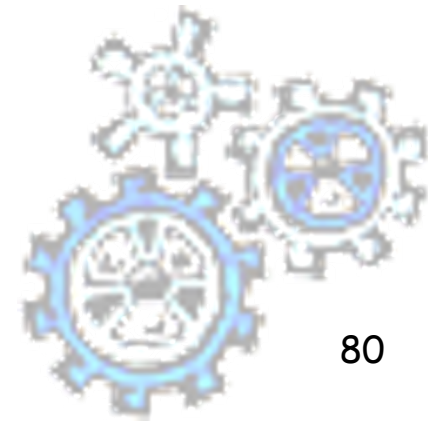
# Mining Association Rules with Constraints

## ■ Postprocessing

- A naïve solution: apply Apriori for finding all frequent sets, and **then** to test them for constraint satisfaction one by one.

## ■ Optimization

- Han approach: comprehensive analysis of the properties of constraints and try to **push them as deeply as possible** inside the frequent set computation.



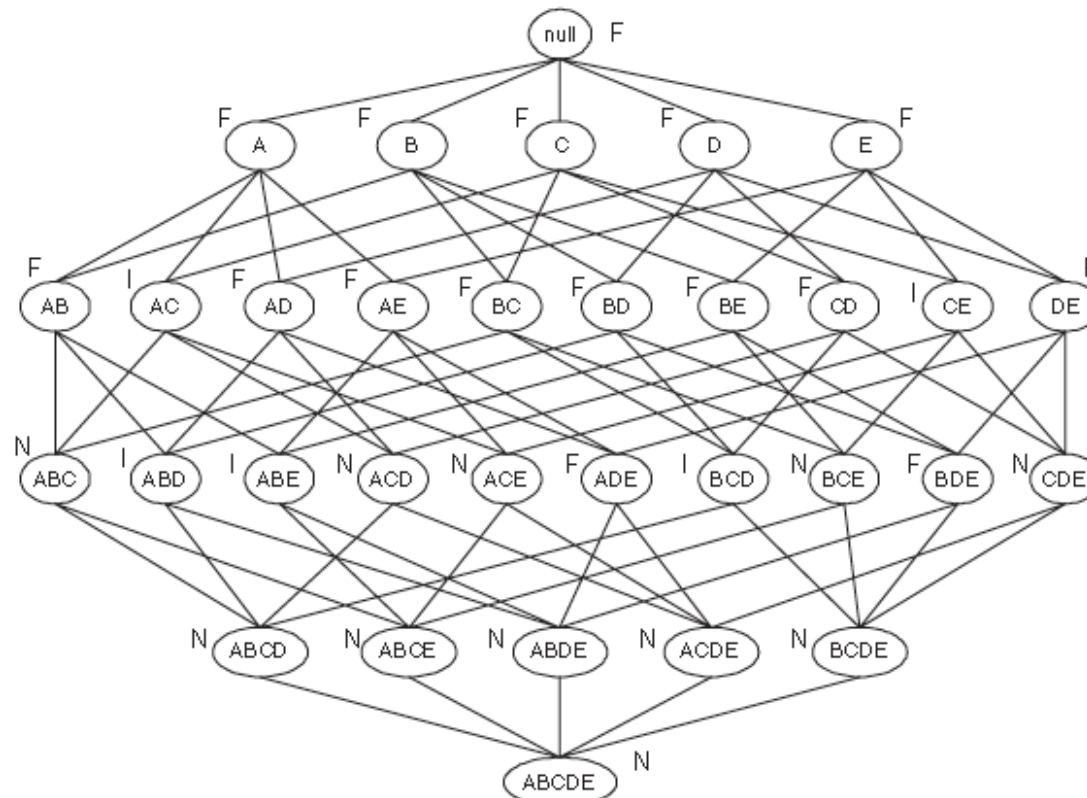
# Exercise 6

TABLE 9.3. Example of market basket transactions.

Transaction ID	Items Bought
1	$\{a, b, d, e\}$
2	$\{b, c, d\}$
3	$\{a, b, d, e\}$
4	$\{a, c, d, e\}$
5	$\{b, c, d, e\}$
6	$\{b, d, e\}$
7	$\{c, d\}$
8	$\{a, b, c\}$
9	$\{a, d, e\}$
10	$\{b, d\}$

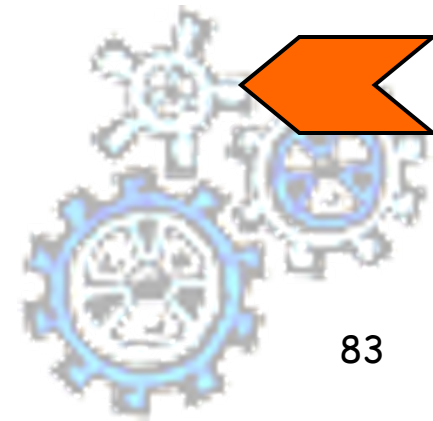


# Exercise 8 Solution



# Association rules - module outline

- What are association rules (AR) and what are they used for:
  - The paradigmatic application: Market Basket Analysis
  - The single dimensional AR (intra-attribute)
- How to compute AR
  - Basic Apriori Algorithm and its optimizations
  - Multi-Dimension AR (inter-attribute)
  - Quantitative AR
  - Constrained AR
- How to reason on AR and how to evaluate their quality
  - Multiple-level AR
  - Interestingness
  - Correlation vs. Association



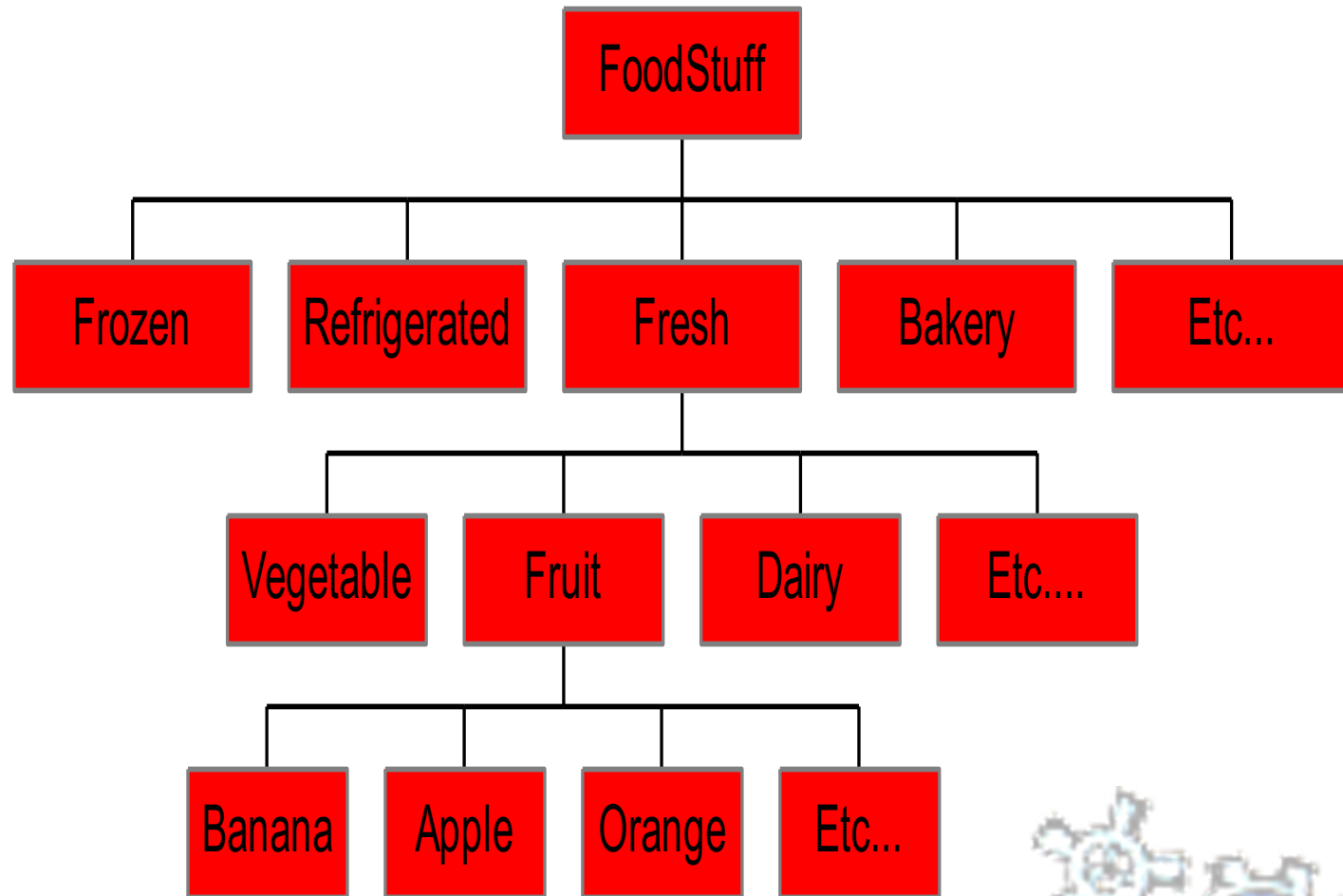
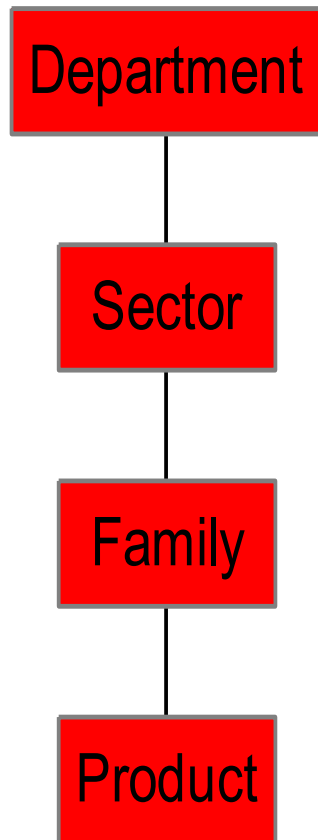


# Multilevel AR

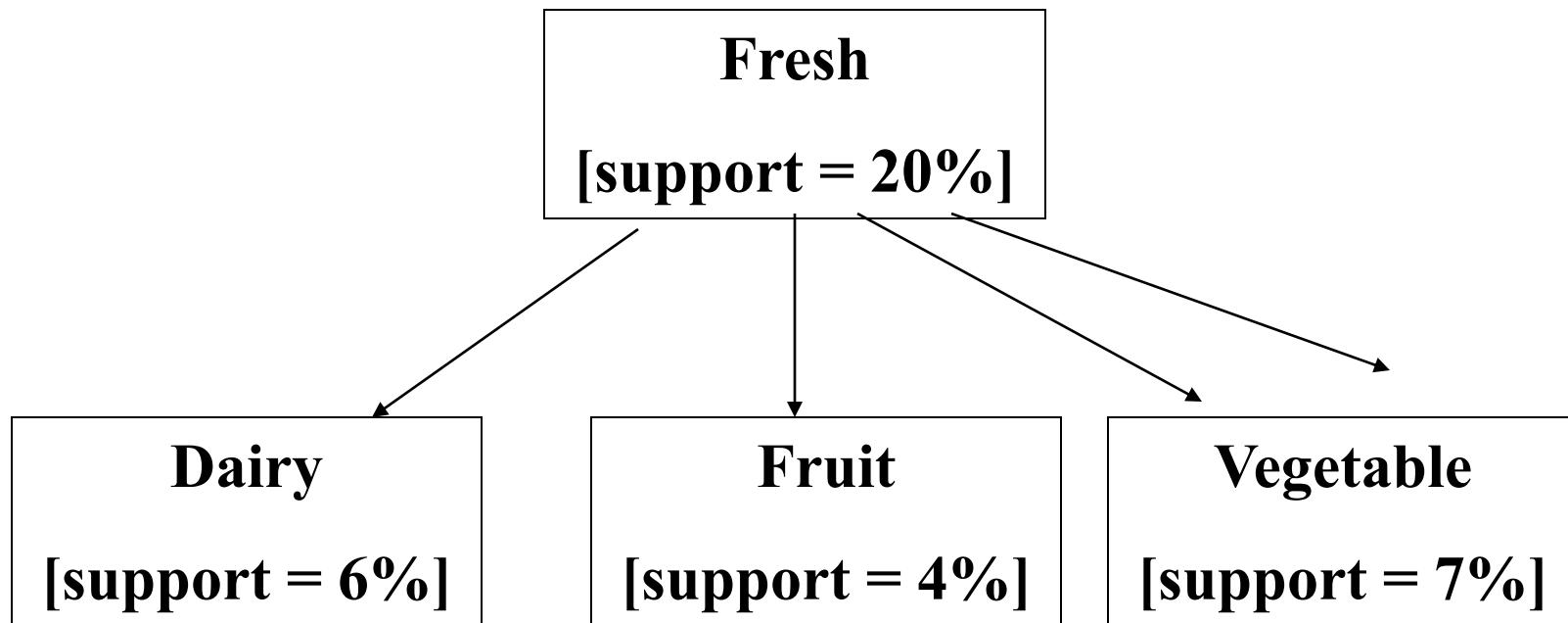
- Is difficult to find interesting patterns at a **too primitive level**
  - high support = too few rules
  - low support = too many rules, most uninteresting
- Approach: reason at suitable level of abstraction
- A common form of background knowledge is that an attribute may be generalized or specialized according to a **hierarchy of concepts**
- Dimensions and levels can be efficiently encoded in transactions
- **Multilevel Association Rules** : rules which combine associations with hierarchy of concepts



# Hierarchy of concepts



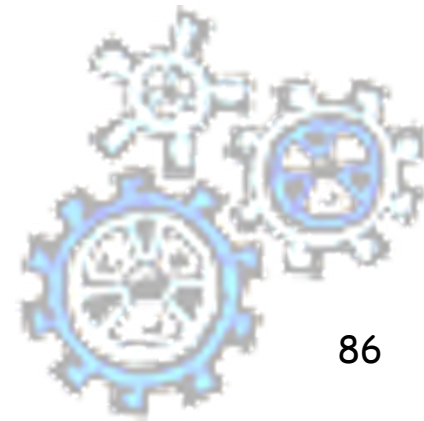
# Multilevel AR



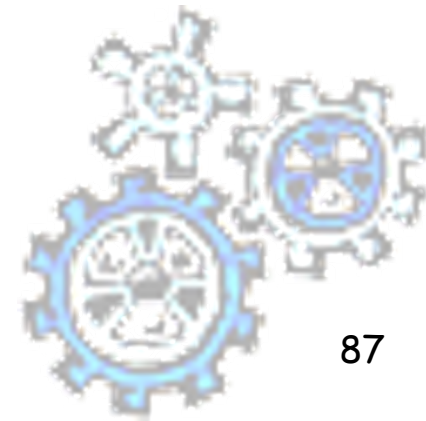
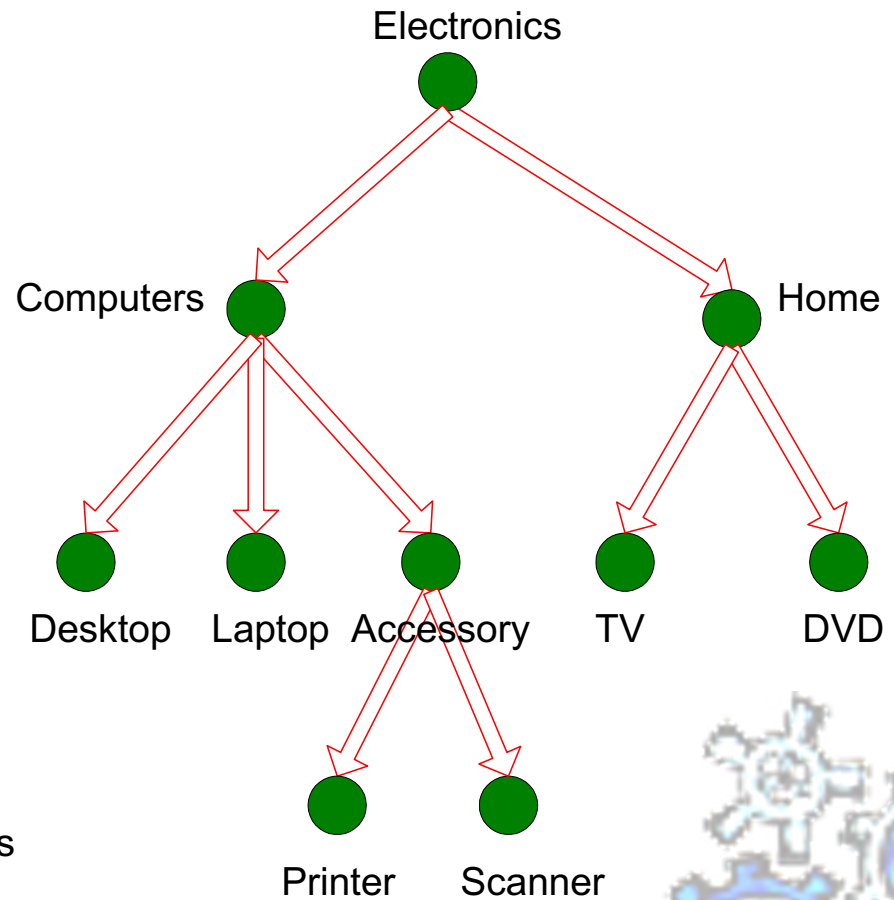
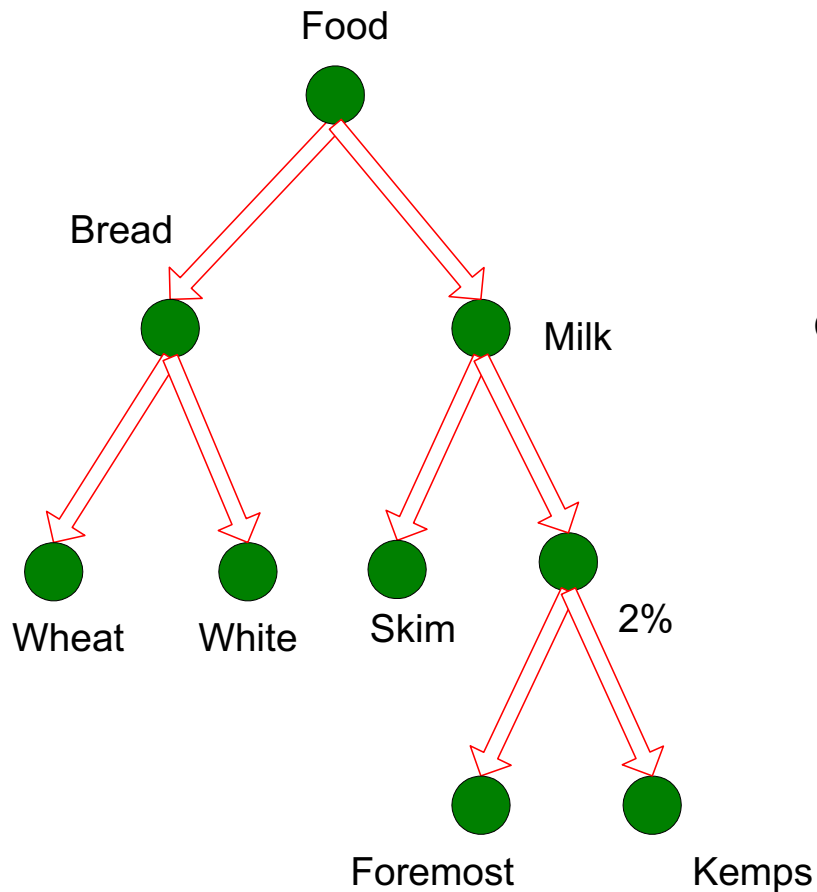
**Fresh  $\Rightarrow$  Bakery [20%, 60%]**

**Dairy  $\Rightarrow$  Bread [6%, 50%]**

**Fruit  $\Rightarrow$  Bread [1%, 50%] is not valid**

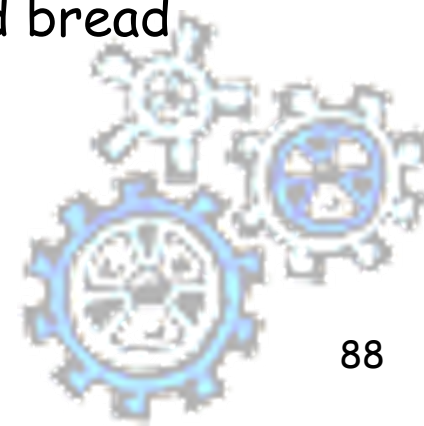


# Multi-level Association Rules



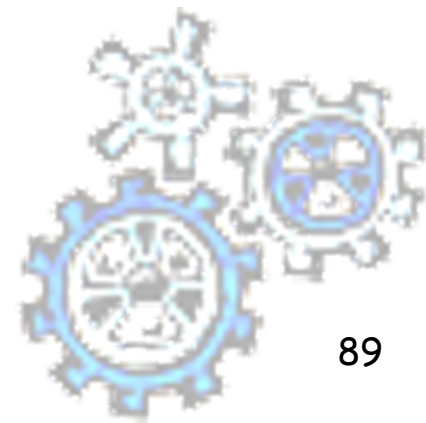
# Multi-level Association Rules

- Why should we incorporate concept hierarchy?
    - Rules at lower levels may not have enough support to appear in any frequent itemsets
    - Rules at lower levels of the hierarchy are overly specific
      - ✓ e.g., skim milk → white bread, 2% milk → wheat bread, skim milk → wheat bread, etc.
- are indicative of association between milk and bread



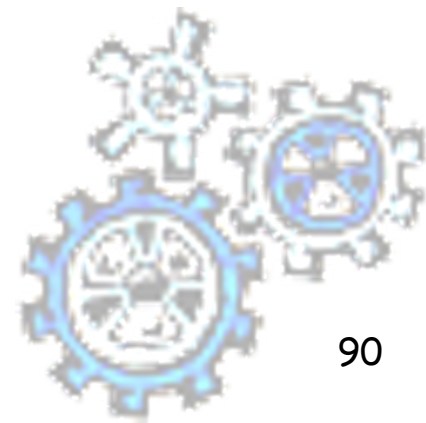
# Support and Confidence of Multilevel AR

- **from specialized to general:** support of rules increases (new rules may become valid)
- **from general to specialized:** support of rules decreases (rules may become not valid, their support falls under the threshold)
- **Confidence is not affected**



# Multi-level Association Rules

- How do support and confidence vary as we traverse the concept hierarchy?
  - If  $X$  is the parent item for both  $X1$  and  $X2$ , then
$$\sigma(X) \leq \sigma(X1) + \sigma(X2)$$
  - If  $\sigma(X1 \cup Y1) \geq \text{minsup}$ ,  
and  $X$  is parent of  $X1$ ,  $Y$  is parent of  $Y1$   
then  $\sigma(X \cup Y1) \geq \text{minsup}$ ,  $\sigma(X1 \cup Y) \geq \text{minsup}$   
 $\sigma(X \cup Y) \geq \text{minsup}$
  - If  $\text{conf}(X1 \Rightarrow Y1) \geq \text{minconf}$ ,  
then  $\text{conf}(X1 \Rightarrow Y) \geq \text{minconf}$



# Reasoning with Multilevel AR

- Too low level => too many rules and too primitive.

Example: **Apple Melinda**  $\Rightarrow$  **Colgate Tooth-paste**

It is a curiosity not a behavior

- Too high level => uninteresting rules

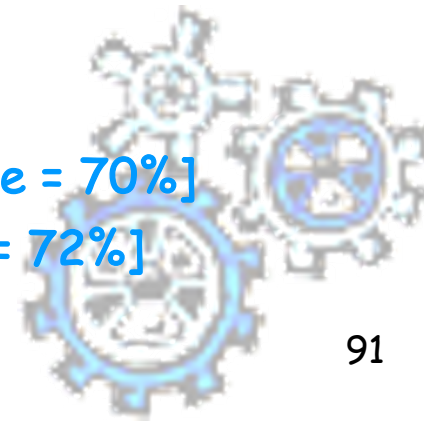
Example: **Foodstuff**  $\Rightarrow$  **Varia**

- Redundancy => some rules may be redundant due to "ancestor" relationships between items.

- A rule is redundant if its support is close to the "expected" value, based on the rule's ancestor.

- Example (milk has 4 subclasses)

- **milk**  $\Rightarrow$  **wheat bread**, [support = 8%, confidence = 70%]
  - **2%-milk**  $\Rightarrow$  **wheat bread**, [support = 2%, confidence = 72%]





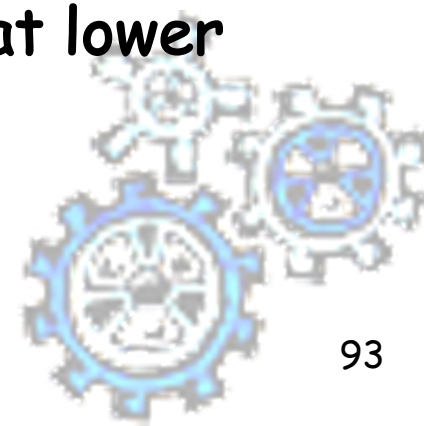
# Mining Multilevel AR

- Calculate frequent itemsets at each concept level, until no more frequent itemsets can be found
- For each level use Apriori
- A top\_down, progressive deepening approach:
  - First find high-level strong rules:  
fresh → bakery [20%, 60%].
  - Then find their lower-level “weaker” rules:  
fruit → bread [6%, 50%].
- Variations at mining multiple-level association rules.
  - Level-crossed association rules:  
fruit → *wheat bread*
  - Association rules with multiple, alternative hierarchies:  
fruit → *Wonder bread*



# Multi-level Association: Uniform Support vs. Reduced Support

- **Uniform Support: the same minimum support for all levels**
  - + One minimum support threshold. No need to examine itemsets containing any item whose ancestors do not have minimum support.
  - - If support threshold
    - too high  $\Rightarrow$  miss low level associations.
    - too low  $\Rightarrow$  generate too many high level associations.
- **Reduced Support: reduced minimum support at lower levels - different strategies possible**

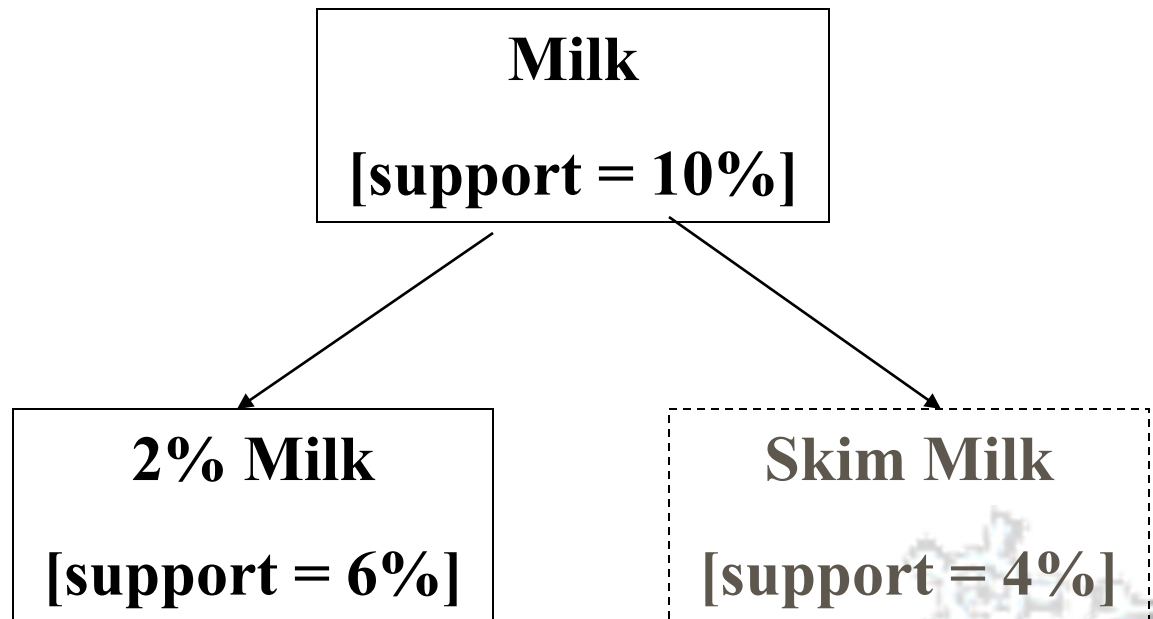


# Uniform Support

## Multi-level mining with uniform support

**Level 1**  
**min\_sup = 5%**

**Level 2**  
**min\_sup = 5%**

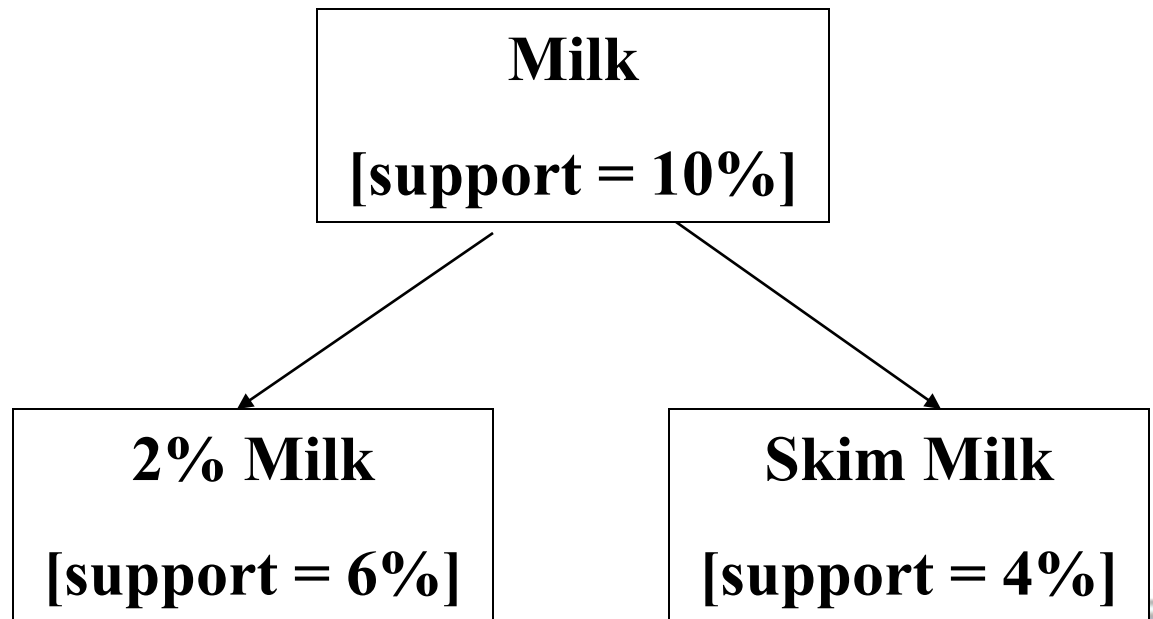


# Reduced Support

## Multi-level mining with reduced support

**Level 1**  
**min\_sup = 5%**

**Level 2**  
**min\_sup = 3%**



# Beyond Support and Confidence

## ■ Example 1: (Aggarwal & Yu, PODS98)

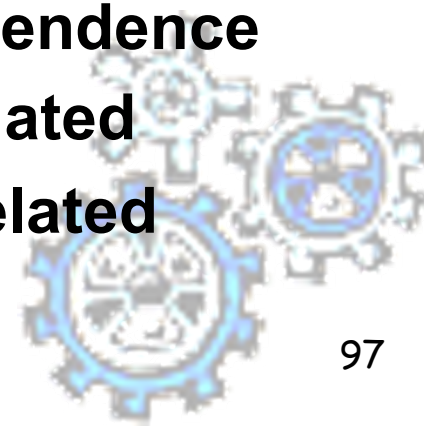
	coffee	not coffee	sum(row)
tea	20	5	25
not tea	70	5	75
sum(col.)	90	10	100

- $\{tea\} \Rightarrow \{coffee\}$  has high support (20%) and confidence (80%)
- However, a priori probability that a customer buys coffee is 90%
  - A customer who is known to buy tea is less likely to buy coffee (by 10%)
  - There is a negative correlation between buying tea and buying coffee
  - $\{\sim tea\} \Rightarrow \{coffee\}$  has higher confidence(93%)



# Statistical Independence

- Population of 1000 students
  - 600 students know how to swim (S)
  - 700 students know how to bike (B)
  - 420 students know how to swim and bike (S,B)
- $P(S \wedge B) = 420/1000 = 0.42$
- $P(S) \times P(B) = 0.6 \times 0.7 = 0.42$
- $P(S \wedge B) = P(S) \times P(B) \Rightarrow$  Statistical independence
- $P(S \wedge B) > P(S) \times P(B) \Rightarrow$  Positively correlated
- $P(S \wedge B) < P(S) \times P(B) \Rightarrow$  Negatively correlated



# Correlation and Interest

- Two events are independent if  $P(A \wedge B) = P(A) * P(B)$ , otherwise are correlated.
- Interest =  $P(A \wedge B) / P(B) * P(A)$
- Interest expresses measure of correlation
  - $= 1 \Rightarrow A$  and  $B$  are independent events
  - **less than 1**  $\Rightarrow A$  and  $B$  negatively correlated,
  - **greater than 1**  $\Rightarrow A$  and  $B$  positively correlated.
  - In our example,  $I(\text{buy tea} \wedge \text{buy coffee}) = 0.89$  i.e. they are negatively correlated.



# Computing Interestingness Measure

- Given a rule  $X \rightarrow Y$ , information needed to compute rule interestingness can be obtained from a contingency table

Contingency table for  $X \rightarrow Y$

	$Y$	$\overline{Y}$	
$X$	$f_{11}$	$f_{10}$	$f_{1+}$
$\overline{X}$	$f_{01}$	$f_{00}$	$f_{0+}$
	$f_{+1}$	$f_{+0}$	$ T $

$f_{11}$ : support of  $X$  and  $Y$

$f_{10}$ : support of  $X$  and  $\overline{Y}$

$f_{01}$ : support of  $\overline{X}$  and  $Y$

$f_{00}$ : support of  $\overline{X}$  and  $\overline{Y}$

Used to define various measures

- ◆ support, confidence, lift, Gini, J-measure, etc.



# Statistical-based Measures

- Measures that take into account statistical dependence

$$Lift = \frac{P(Y | X)}{P(Y)}$$

$$Interest = \frac{P(X, Y)}{P(X)P(Y)}$$

$$PS = P(X, Y) - P(X)P(Y)$$

$$\phi - coefficient = \frac{P(X, Y) - P(X)P(Y)}{\sqrt{P(X)[1 - P(X)]P(Y)[1 - P(Y)]}}$$

# Example: Lift/Interest

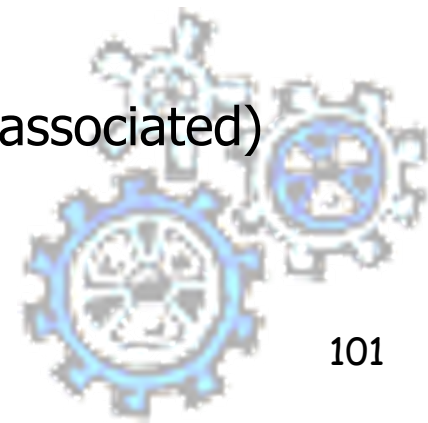
	Coffee	<u>Coffee</u>	
<del>Tea</del>	15	5	20
Tea	75	5	80
	90	10	100

Association Rule: Tea  $\rightarrow$  Coffee

Confidence=  $P(\text{Coffee}|\text{Tea}) = 0.75$

but  $P(\text{Coffee}) = 0.9$

$\Rightarrow$  Lift =  $0.75/0.9 = 0.8333$  ( $< 1$ , therefore is negatively associated)



# Drawback of Lift & Interest

	y	$\bar{y}$	
X	10	0	10
$\bar{X}$	0	90	90
	10	90	100

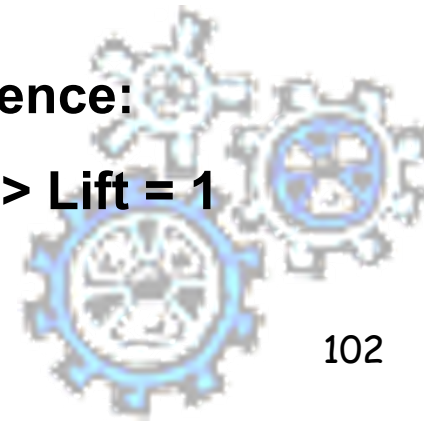
	y	$\bar{y}$	
X	90	0	90
$\bar{X}$	0	10	10
	90	10	100

$$Lift = \frac{0.1}{(0.1)(0.1)} = 10$$

$$Lift = \frac{0.9}{(0.9)(0.9)} = 1.11$$

**Statistical independence:**

**If  $P(X,Y)=P(X)P(Y) \Rightarrow Lift = 1$**



There are lots of measures proposed in the literature

Some measures are good for certain applications, but not for others

What criteria should we use to determine whether a measure is good or bad?

What about Apriori-style support based pruning? How does it affect these measures?

#	Measure	Formula
1	$\phi$ -coefficient	$\frac{P(A,B) - P(A)P(B)}{\sqrt{P(A)P(B)(1-P(A))(1-P(B))}}$
2	Goodman-Kruskal's ( $\lambda$ )	$\frac{\sum_j \max_k P(A_j, B_k) + \sum_k \max_j P(A_j, B_k) - \max_j P(A_j) - \max_k P(B_k)}{2 - \max_j P(A_j) - \max_k P(B_k)}$
3	Odds ratio ( $\alpha$ )	$\frac{P(A,B)P(\bar{A},\bar{B})}{P(A,\bar{B})P(\bar{A},B)}$
4	Yule's $Q$	$\frac{P(A,B)P(\bar{A}\bar{B}) - P(A,\bar{B})P(\bar{A},B)}{P(A,B)P(\bar{A}\bar{B}) + P(A,\bar{B})P(\bar{A},B)} = \frac{\alpha - 1}{\alpha + 1}$
5	Yule's $Y$	$\frac{\sqrt{P(A,B)P(\bar{A}\bar{B})} - \sqrt{P(A,\bar{B})P(\bar{A},B)}}{\sqrt{P(A,B)P(\bar{A}\bar{B})} + \sqrt{P(A,\bar{B})P(\bar{A},B)}} = \frac{\sqrt{\alpha} - 1}{\sqrt{\alpha} + 1}$
6	Kappa ( $\kappa$ )	$\frac{P(A,B) + P(\bar{A},\bar{B}) - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}$
7	Mutual Information ( $M$ )	$\frac{\sum_i \sum_j P(A_i, B_j) \log \frac{P(A_i, B_j)}{P(A_i)P(B_j)}}{\min(-\sum_i P(A_i) \log P(A_i), -\sum_j P(B_j) \log P(B_j))}$
8	J-Measure ( $J$ )	$\max \left( P(A, B) \log \left( \frac{P(B A)}{P(B)} \right) + P(\bar{A}\bar{B}) \log \left( \frac{P(\bar{B} \bar{A})}{P(\bar{B})} \right), \right. \\ \left. P(A, B) \log \left( \frac{P(A B)}{P(A)} \right) + P(\bar{A}\bar{B}) \log \left( \frac{P(\bar{A} \bar{B})}{P(\bar{A})} \right) \right)$
9	Gini index ( $G$ )	$\max \left( P(A)[P(B A)^2 + P(\bar{B} A)^2] + P(\bar{A})[P(B \bar{A})^2 + P(\bar{B} \bar{A})^2] \right. \\ \left. - P(B)^2 - P(\bar{B})^2, \right. \\ \left. P(B)[P(A B)^2 + P(\bar{A} B)^2] + P(\bar{B})[P(A \bar{B})^2 + P(\bar{A} \bar{B})^2] \right. \\ \left. - P(A)^2 - P(\bar{A})^2 \right)$
10	Support ( $s$ )	$P(A, B)$
11	Confidence ( $c$ )	$\max(P(B A), P(A B))$
12	Laplace ( $L$ )	$\max \left( \frac{NP(A,B)+1}{NP(A)+2}, \frac{NP(A,B)+1}{NP(B)+2} \right)$
13	Conviction ( $V$ )	$\max \left( \frac{P(A)P(\bar{B})}{P(\bar{A}B)}, \frac{P(B)P(\bar{A})}{P(\bar{B}A)} \right)$
14	Interest ( $I$ )	$\frac{P(A,B)}{P(A)P(B)}$
15	cosine ( $IS$ )	$\frac{P(A,B)}{\sqrt{P(A)P(B)}}$
16	Piatetsky-Shapiro's ( $PS$ )	$P(A, B) - P(A)P(B)$
17	Certainty factor ( $F$ )	$\max \left( \frac{P(B A) - P(B)}{1 - P(B)}, \frac{P(A B) - P(A)}{1 - P(A)} \right)$
18	Added Value ( $AV$ )	$\max(P(B A) - P(B), P(A B) - P(A))$
19	Collective strength ( $S$ )	$\frac{P(A,B) + P(\bar{A}\bar{B})}{P(A)P(B) + P(\bar{A})P(\bar{B})} \times \frac{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A,B) - P(\bar{A}\bar{B})}$
20	Jaccard ( $\zeta$ )	$\frac{P(A,B)}{P(A) + P(B) - P(A,B)}$
21	Klosgen ( $K$ )	$\sqrt{P(A, B)} \max(P(B A) - P(B), P(A B) - P(A))$

# Properties of A Good Measure

## ■ Piatetsky-Shapiro:

3 properties a good measure  $M$  must satisfy:

- $M(A,B) = 0$  if  $A$  and  $B$  are statistically independent
- $M(A,B)$  increase monotonically with  $P(A,B)$  when  $P(A)$  and  $P(B)$  remain unchanged
- $M(A,B)$  decreases monotonically with  $P(A)$  [or  $P(B)$ ] when  $P(A,B)$  and  $P(B)$  [or  $P(A)$ ] remain unchanged



# Comparing Different Measures

10 examples of  
contingency tables:

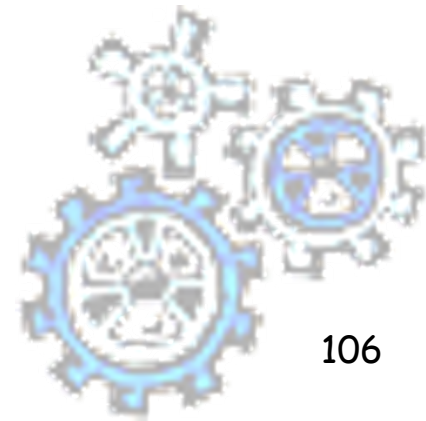
Example	$f_{11}$	$f_{10}$	$f_{01}$	$f_{00}$
E1	8123	83	424	1370
E2	8330	2	622	1046
E3	9481	94	127	298
E4	3954	3080	5	2961
E5	2886	1363	1320	4431
E6	1500	2000	500	6000
E7	4000	2000	1000	3000
E8	4000	2000	2000	2000
E9	1720	7121	5	1154
E10	61	2483	4	7452

Rankings of contingency tables  
using various measures:

#	$\phi$	$\lambda$	$\alpha$	$Q$	$Y$	$\kappa$	$M$	$J$	$G$	$s$	$c$	$L$	$V$	$I$	$IS$	$PS$	$F$	$AV$	$S$	$\zeta$	$K$
E1	1	1	3	3	3	1	2	2	1	3	5	5	4	6	2	2	4	6	1	2	5
E2	2	2	1	1	1	2	1	3	2	2	1	1	1	8	3	5	1	8	2	3	6
E3	3	3	4	4	4	3	3	8	7	1	4	4	6	10	1	8	6	10	3	1	10
E4	4	7	2	2	2	5	4	1	3	6	2	2	2	4	4	1	2	3	4	5	1
E5	5	4	8	8	8	4	7	5	4	7	9	9	9	3	6	3	9	4	5	6	3
E6	6	6	7	7	7	7	6	4	6	9	8	8	7	2	8	6	7	2	7	8	2
E7	7	5	9	9	9	6	8	6	5	4	7	7	8	5	5	4	8	5	6	4	4
E8	8	9	10	10	10	8	10	10	8	4	10	10	10	9	7	7	10	9	8	7	9
E9	9	9	5	5	5	9	9	7	9	8	3	3	3	7	9	9	3	7	9	9	8
E10	10	8	6	6	6	10	5	9	10	10	6	6	5	1	10	10	5	1	10	10	7

# Domain dependent measures

- Together with support, confidence, interest, ..., use also (in post-processing) domain-dependent measures
- E.g., use rule constraints on rules
- Example: take only rules which are significant with respect their economic value
- $\text{sum(LHS)} + \text{sum(RHS)} > 100$



# MBA in Web Usage Mining

## ■ Association Rules in Web Transactions

- discover affinities among sets of Web page references across user sessions

## ■ Examples

- 60% of clients who accessed `/products/`, also accessed `/products/software/webminer.htm`
- 30% of clients who accessed `/special-offer.html`, placed an online order in `/products/software/`
- Actual Example from IBM official Olympics Site:
  - ✓ {Badminton, Diving} ==> {Table Tennis}
  - [conf = 69.7%, sup = 0.35%]

## ■ Applications

- Use rules to serve dynamic, customized contents to users
- prefetch files that are most likely to be accessed
- determine the best way to structure the Web site (site optimization)

- targeted electronic advertising and increasing cross sales



# Web Usage Mining: Example

## ■ Association Rules From Cray Research Web Site

Conf	supp	Association Rule
82.8	3.17	/PUBLIC/product-info/T3E ====> /PUBLIC/product-info/T3E/CRAY_T3E.html
90	0.14	/PUBLIC/product-info/J90/J90.html, /PUBLIC/product-info/T3E ====> /PUBLIC/product-info/T3E/CRAY_T3E.html
97.2	0.15	/PUBLIC/product-info/J90, /PUBLIC/product-info/T3E/CRAY_T3E.html, /PUBLIC/product-info/T90, ====> /PUBLIC/product-info/T3E, /PUBLIC/sc.html

## ■ Design "suggestions"

- from rules 1 and 2: there is something in **J90.html** that should be moved to the page **/PUBLIC/product-info/T3E** (why?)



# MBA in Text / Web Content Mining

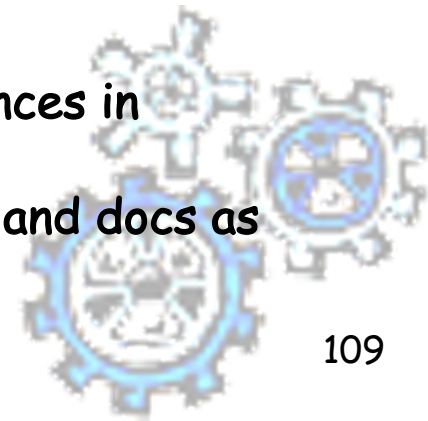
## ■ Documents Associations

- Find (content-based) associations among documents in a collection
- Documents correspond to items and words correspond to transactions
- Frequent itemsets are groups of docs in which many words occur in common

	Doc 1	Doc 2	Doc 3	...	Doc n
business	5	5	2	...	1
capital	2	4	3	...	5
fund	0	0	0	...	1
⋮	⋮	⋮	⋮	⋮	⋮
invest	6	0	0	...	3

## ■ Term Associations

- Find associations among words based on their occurrences in documents
- similar to above, but invert the table (terms as items, and docs as transactions)



# Atherosclerosis prevention study

**2nd Department of Medicine, 1st Faculty of  
Medicine of Charles University and Charles  
University Hospital, U nemocnice 2, Prague  
2 (head. Prof. M. Aschermann, MD, SDr,  
FESC)**

# Atherosclerosis prevention study:

- The STULONG 1 data set is a real database that keeps information about the study of the development of atherosclerosis risk factors in a population of middle aged men.
- Used for Discovery Challenge at PKDD 00-02-03-04



# Atherosclerosis prevention study:

- Study on 1400 middle-aged men at Czech hospitals
  - Measurements concern development of cardiovascular disease and other health data in a series of exams
- The aim of this analysis is to look for associations between medical characteristics of patients and death causes.
- Four tables
  - Entry and subsequent exams, questionnaire responses, deaths



# The input data

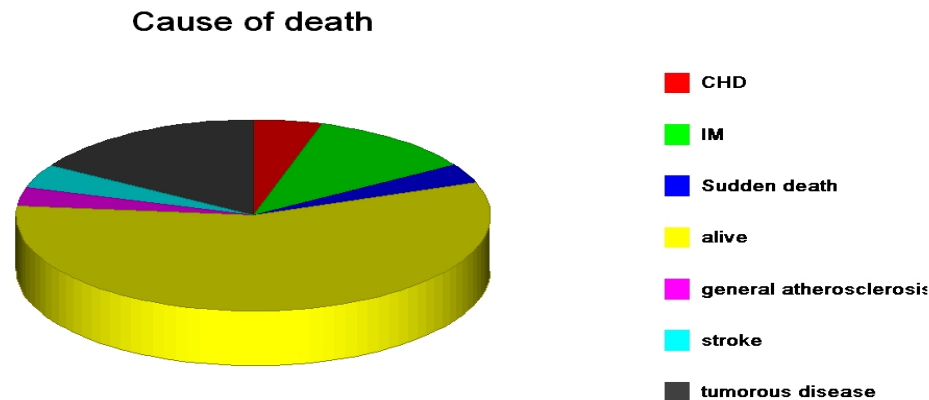
Data from Entry and Exams		
General characteristics	Examinations	habits
Marital status	Chest pain	Alcohol
Transport to a job	Breathlessness	Liquors
Physical activity in a job	Cholesterol	Beer 10
Activity after a job	Urine	Beer 12
Education	Subscapular	Wine
Responsibility	Triceps	Smoking
Age		Former smoker
Weight		Duration of smoking
Height		Tea
		Sugar
		Coffee

# The input data

DEATH CAUSE	PATIENTS	%
myocardial infarction	80	20.6
coronary heart disease	33	8.5
stroke	30	7.7
other causes	79	20.3
sudden death	23	5.9
unknown	8	2.0
tumorous disease	114	29.3
general atherosclerosis	22	5.7
TOTAL	389	100.0

# Data selection

- When joining “Entry” and “Death” tables we implicitly create a new attribute “Cause of death”, which is set to “alive” for subjects present in the “Entry” table but not in the “Death” table.
- We have only 389 subjects in death table.





# The prepared data

Patient	General characteristics		Examinations		Habits		Cause of death
	Activity after work	Education	Chest pain	...	Alcohol	.....	
1	moderate activity	university	not present		no		Stroke
2	great activity		not ischaemic		occasionally		myocardial infarction
3	he mainly sits		other pains		regularly		tumorous disease
.....	.....	.....	.....	..	...	.....	alive
389	he mainly sits		other pains		regularly		tumorous disease

# Descriptive Analysis/ Subgroup Discovery / Association Rules

Are there strong relations concerning death cause?

General characteristics (?)  $\Rightarrow$  Death cause (?)

Examinations (?)  $\Rightarrow$  Death cause (?)

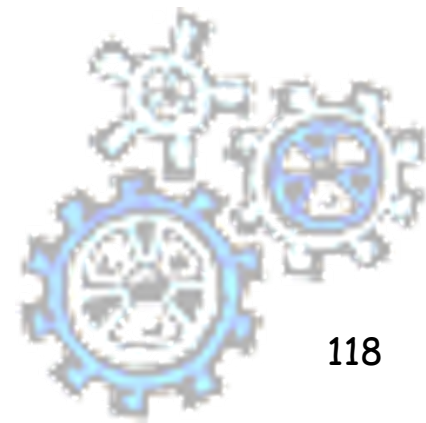
Habits (?)  $\Rightarrow$  Death cause (?)

Combinations (?)  $\Rightarrow$  Death cause (?)



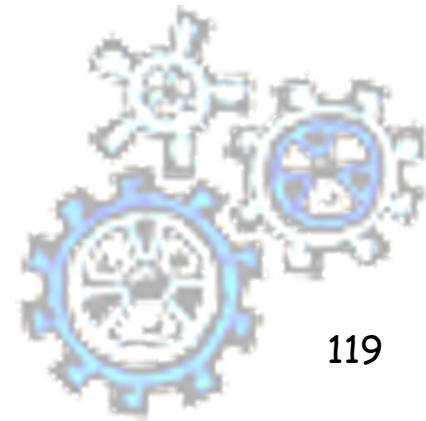
# Example of extracted rules

- **Education(university) & Height<176-180>  
⇒Death cause (tumouros disease), 16 ; 0.62**
- **It means that on tumorous disease have died 16,  
i.e. 62% of patients with university education and  
with height 176-180 cm.**



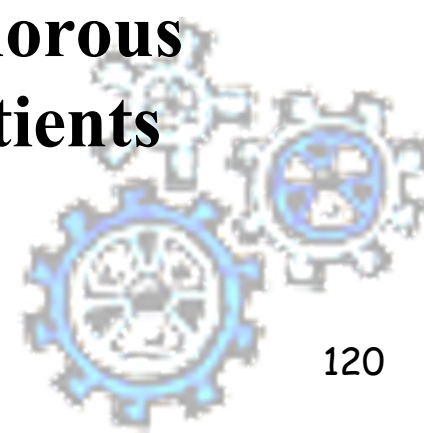
# Example of extracted rules

- **Physical activity in work(he mainly sits) & Height<176-180>  $\Rightarrow$  Death cause (tumouros disease), 24; 0.52**
- **It means that on tumorous disease have died 24 i.e. 52% of patients that mainly sit in the work and whose height is 176-180 cm.**



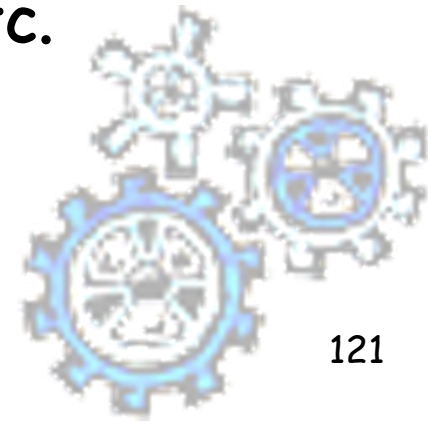
# Example of extracted rules

- **Education(university) & Height<176-180>  
⇒Death cause (tumouros disease),  
*16; 0.62; +1.1;***
- **the relative frequency of patients who died on tumorous disease among patients with university education and with height 176-180 cm is 110 per cent higher than the relative frequency of patients who died on tumorous disease among all the 389 observed patients**



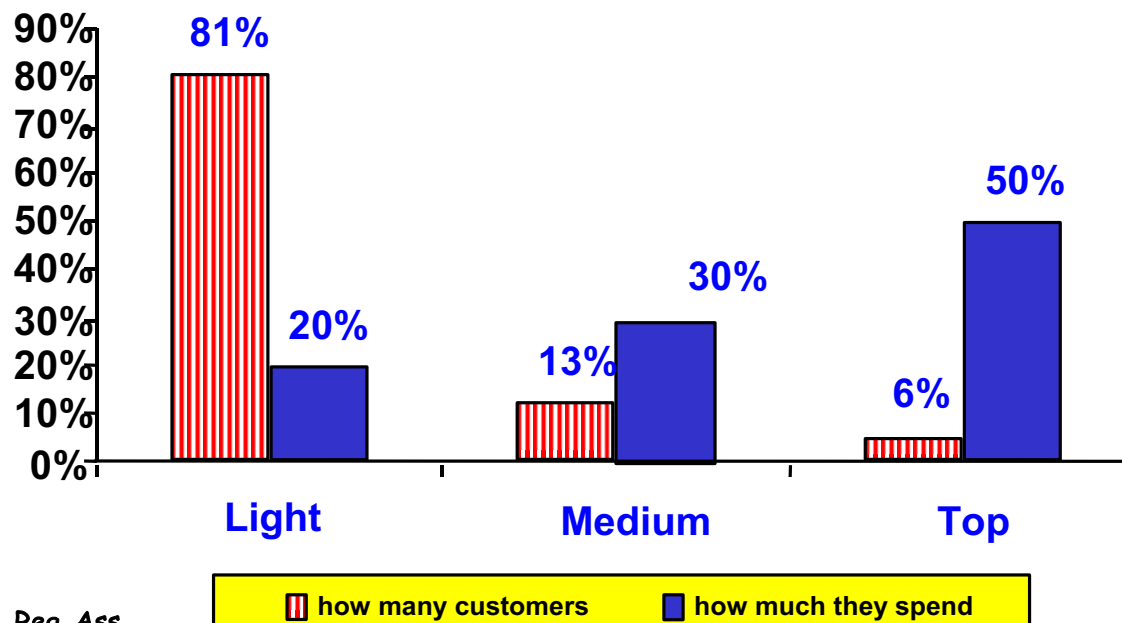
# Conclusions

- **Association rule mining**
  - probably the most significant contribution from the database community to KDD
  - A large number of papers have been published
- **Many interesting issues have been explored**
- **An interesting research direction**
  - Association analysis in other types of data: spatial data, multimedia data, time series data, etc.



# Conclusion (2)

- MBA is a key factor of success in the competition of supermarket retailers.
- Knowledge of customers and their purchasing behavior brings potentially huge added value.



# Which tools for market basket analysis?

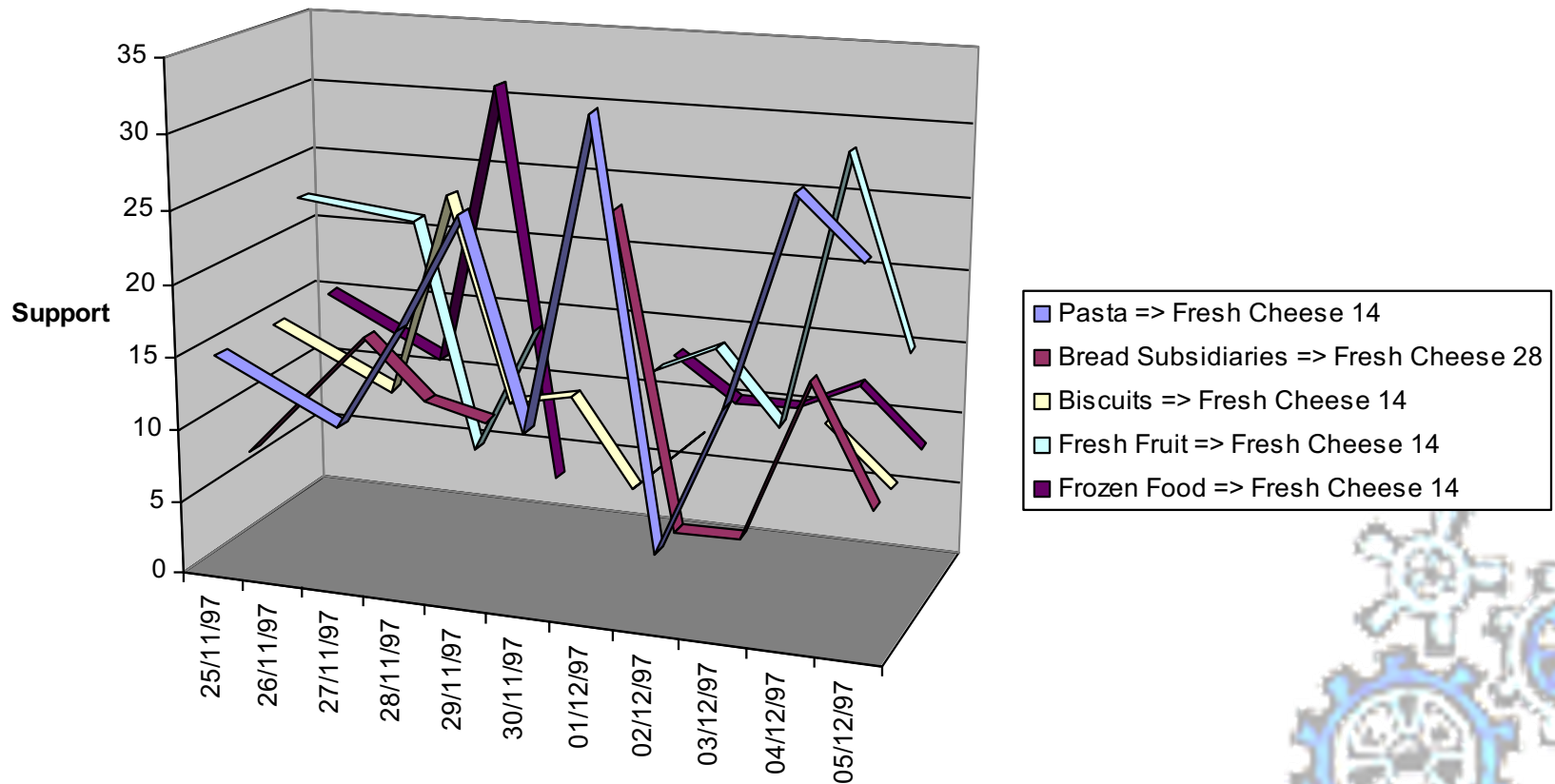
- Association rule are needed but insufficient
- Market analysts ask for **business rules**:
  - Is supermarket assortment adequate for the company's target class of customers?
  - Is a promotional campaign effective in establishing a desired purchasing habit?





# Business rules: temporal reasoning on AR

- Which rules are established by a promotion?
- How do rules change along time?



# References - Association rules

- R. Agrawal, T. Imielinski, and A. Swami. Mining association rules between sets of items in large databases. SIGMOD'93, 207-216, Washington, D.C.
- R. Agrawal and R. Srikant. Fast algorithms for mining association rules. VLDB'94 487-499, Santiago, Chile.
- R. Agrawal and R. Srikant. Mining sequential patterns. ICDE'95, 3-14, Taipei, Taiwan.
- R. J. Bayardo. Efficiently mining long patterns from databases. SIGMOD'98, 85-93, Seattle, Washington.
- S. Brin, R. Motwani, and C. Silverstein. Beyond market basket: Generalizing association rules to correlations. SIGMOD'97, 265-276, Tucson, Arizona..
- D.W. Cheung, J. Han, V. Ng, and C.Y. Wong. Maintenance of discovered association rules in large databases: An incremental updating technique. ICDE'96, 106-114, New Orleans, LA..
- T. Fukuda, Y. Morimoto, S. Morishita, and T. Tokuyama. Data mining using two-dimensional optimized association rules: Scheme, algorithms, and visualization. SIGMOD'96, 13-23, Montreal, Canada.
- E.-H. Han, G. Karypis, and V. Kumar. Scalable parallel data mining for association rules. SIGMOD'97, 277-288, Tucson, Arizona.
- J. Han and Y. Fu. Discovery of multiple-level association rules from large databases. VLDB'95, 420-431, Zurich, Switzerland.
- M. Kamber, J. Han, and J. Y. Chiang. Metarule-guided mining of multi-dimensional association rules using data cubes. KDD'97, 207-210, Newport Beach, California.
- M. Klemettinen, H. Mannila, P. Ronkainen, H. Toivonen, and A.I. Verkamo. Finding interesting rules from large sets of discovered association rules. CIKM'94, 401-408, Gaithersburg, Maryland.
- R. Ng, L. V. S. Lakshmanan, J. Han, and A. Pang. Exploratory mining and pruning optimizations of constrained associations rules. SIGMOD'98, 13-24, Seattle, Washington.
- B. Ozden, S. Ramaswamy, and A. Silberschatz. Cyclic association rules. ICDE'98, 412-421, Orlando, FL.
- J.S. Park, M.S. Chen, and P.S. Yu. An effective hash-based algorithm for mining association rules. SIGMOD'95, 175-186, San Jose, CA.
- S. Ramaswamy, S. Mahajan, and A. Silberschatz. On the discovery of interesting patterns in association rules. VLDB'98, 368-379, New York, NY.
- S. Sarawagi, S. Thomas, and R. Agrawal. Integrating association rule mining with relational database systems: Alternatives and implications. SIGMOD'98, 343-354, Seattle, WA.



# References - Association rules

- A. Savasere, E. Omiecinski, and S. Navathe. An efficient algorithm for mining association rules in large databases. VLDB'95, 432-443, Zurich, Switzerland.
- C. Silverstein, S. Brin, R. Motwani, and J. Ullman. Scalable techniques for mining causal structures. VLDB'98, 594-605, New York, NY.
- R. Srikant and R. Agrawal. Mining generalized association rules. VLDB'95, 407-419, Zurich, Switzerland.
- R. Srikant and R. Agrawal. Mining quantitative association rules in large relational tables. SIGMOD'96, 1-12, Montreal, Canada.
- R. Srikant, Q. Vu, and R. Agrawal. Mining association rules with item constraints. KDD'97, 67-73, Newport Beach, California.
- D. Tsur, J. D. Ullman, S. Abitboul, C. Clifton, R. Motwani, and S. Nestorov. Query flocks: A generalization of association-rule mining. SIGMOD'98, 1-12, Seattle, Washington.
- B. Ozden, S. Ramaswamy, and A. Silberschatz. Cyclic association rules. ICDE'98, 412-421, Orlando, FL.
- R.J. Miller and Y. Yang. Association rules over interval data. SIGMOD'97, 452-461, Tucson, Arizona.
- J. Han, G. Dong, and Y. Yin. Efficient mining of partial periodic patterns in time series database. ICDE'99, Sydney, Australia.
- F. Giannotti, G. Manco, D. Pedreschi and F. Turini. Experiences with a logic-based knowledge discovery support environment. In Proc. 1999 ACM SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery (SIGMOD'99 DMKD). Philadelphia, May 1999.
- F. Giannotti, M. Nanni, G. Manco, D. Pedreschi and F. Turini. Integration of Deduction and Induction for Mining Supermarket Sales Data. In Proc. PADD'99, Practical Application of Data Discovery, Int. Conference, London, April 1999.

